# AI SUMMIT

COLUMBUS, OH • OCTOBER 25–27, 2022

# BUILDING TRUST IN AI SYSTEMS: WHERE ARE WE NOW?

## A DISCUSSION OF APPROACHES IN MOTION

AFDO RAPS HEALTHCARE PRODUCTS COLLABORATIVE | AFDO RAPS

**AI Summit**
COLUMBUS, OH • OCTOBER 25–27, 2022

**Bakul Patel,** Senior Director, Global Digital Health Strategy & Regulatory, Google

**Koen Cobbaert,** Senior Manager – Quality, Standards & Regulations, Philips

**Rohit Nayak,** Co-Founder CEO, Band Connect Inc. | CEO, Electronic Registry Systems, Inc.

# Discussion Outline

- **Why Trust?**
  - Why is it important?
  - Different names and flavors  - Transparency, Trustworthiness, Explainability
- **Driving Forces**
  - GDPR, XAI, EU AI ACT , AI Bill of Rights

This discussion reviews both the regulatory policy as well as the steps being taken by a significant industry player

# AI Bill of Rights

**Safe and Effective Systems**

**Algorithmic Discrimination Protections**

**Data Privacy**

**Notice and Explanation**

**Human Alternatives, Consideration, and Fallback**

https://www.whitehouse.gov/ostp/ai-bill-of-rights/

**BLUEPRINT FOR AN AI BILL OF RIGHTS**

**MAKING AUTOMATED SYSTEMS WORK FOR THE AMERICAN PEOPLE**

OCTOBER 2022

**FROM PRINCIPLES TO PRACTICE**

**A TECHINCAL COMPANION TO THE BLUEPRINT FOR AN AI BILL OF RIGHTS**

# Explainable AI

- EU GDPR – Right of Explanation
- IEEE - Standard for XAI – eXplainable Artificial Intelligence - for Achieving Clarity and Interoperability of AI Systems Design
- IEEE - Guide for an Architectural Framework for Explainable Artificial Intelligence
- DARPA's XAI Initiative
- …..

# EU AI ACT
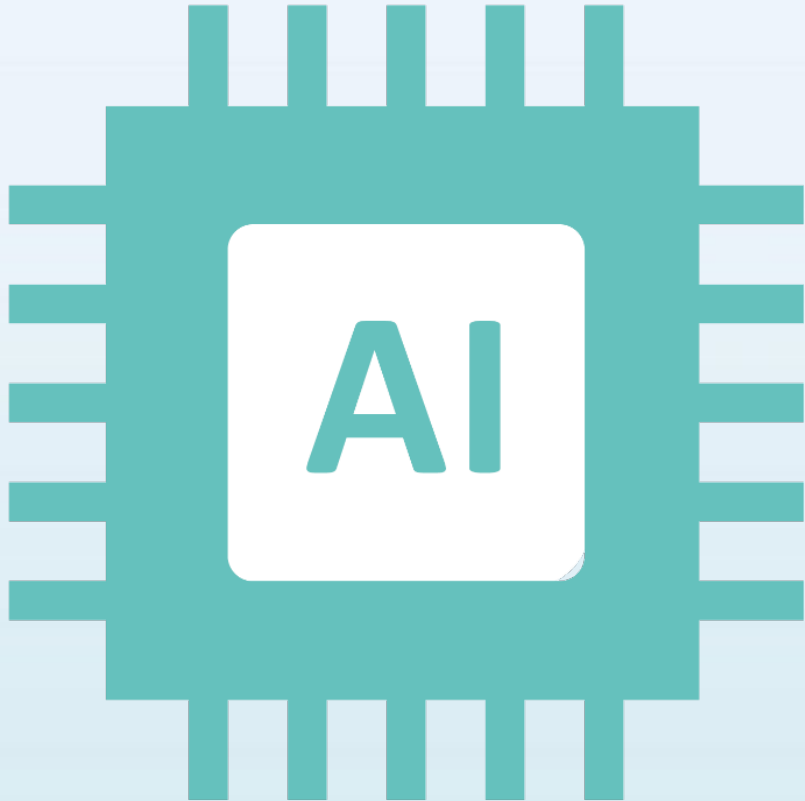
**What is it?**

**Why are we talking about it?**

**Where does it stand?**

legislative draft available
currently in legislative process
adoption expected Q2 2023
data of application Q2 2025 (?)

# AI System – Commission Definition

Proposed AIA Art. 3(1)

**an AI system is**

software that is developed with one or more of the techniques and approaches listed in Annex I and can, for a given set of human-defined objectives, generate outputs such as content, predictions, recommendations, or decisions influencing the environments they interact with.

**Annex I (can be updated through delegated act)**

a) Machine learning approaches, including supervised, unsupervised and reinforcement learning, using a wide variety of methods including deep learning;

b) Logic- and knowledge-based approaches, including knowledge representation, inductive (logic) programming, knowledge bases, inference and deductive engines, (symbolic) reasoning and expert systems;

c) Statistical approaches, Bayesian estimation, search and optimization methods

**definition**

*reads as*

# AI system = any software application

In listing technologies considered AI, Annex I tries to compensate for a vague AI system definition,
but as technologies can be added or removed over time, it increases legal uncertainty
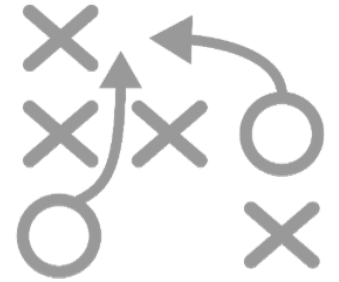


machine learning



inference and deductive engines



reasoning and expert systems



search & optimization methods



logic- and inductive programming
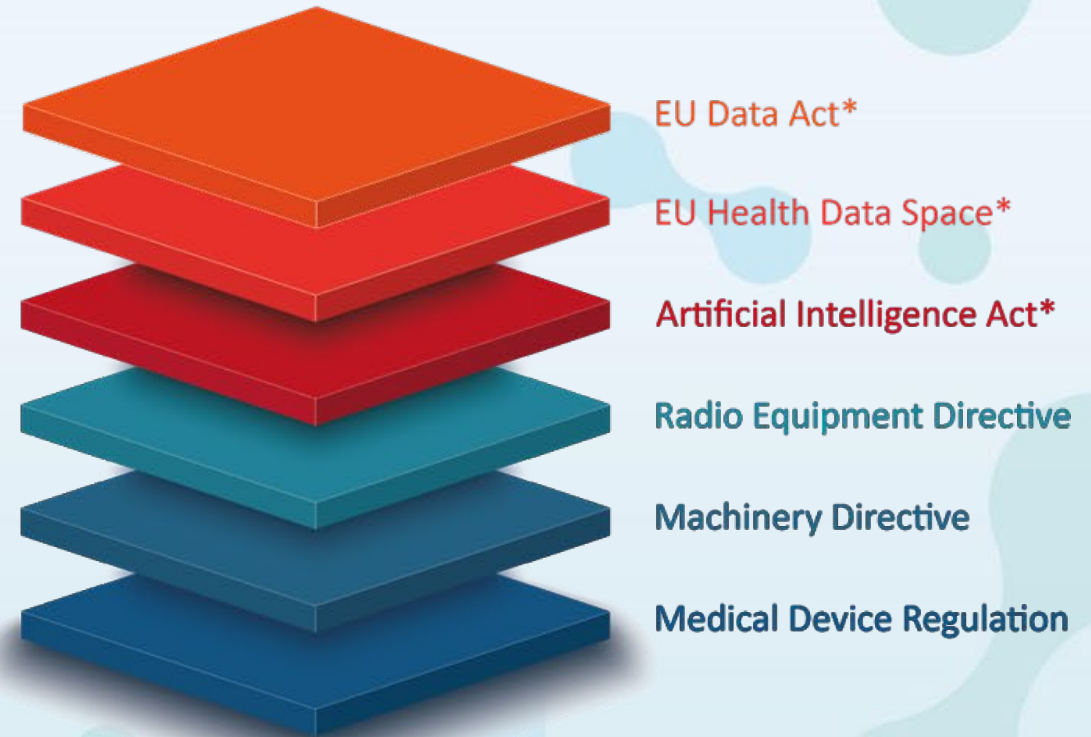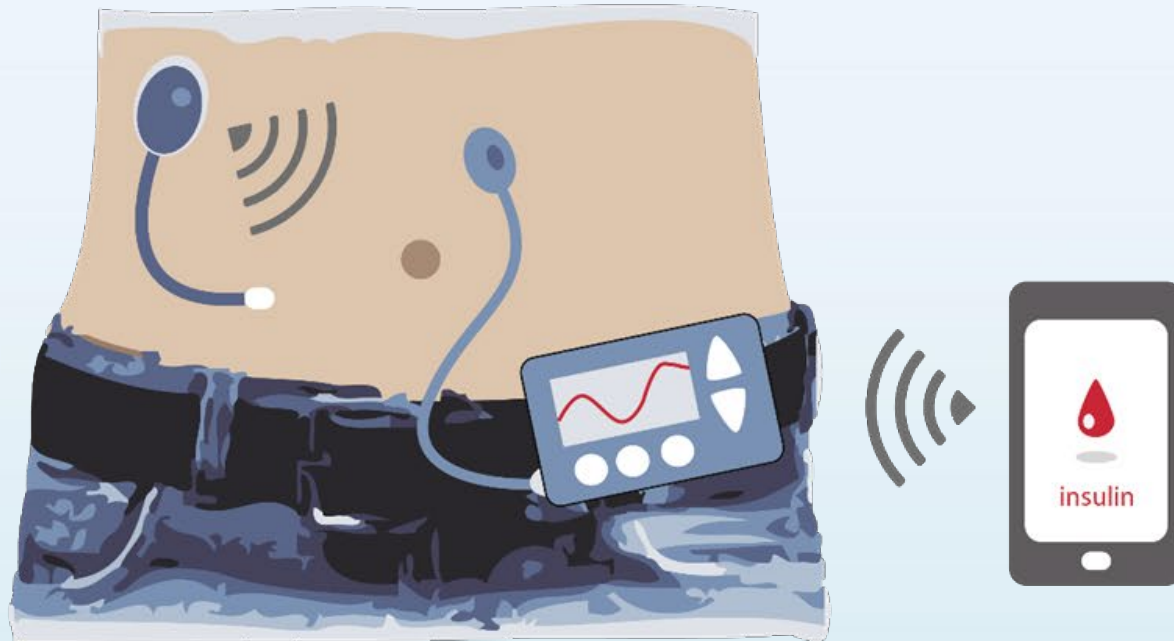
**AI Act contains mandatory requirements for**
**High-Risk AI systems**
**=**
regulated products or safety components of regulated products
which are subject to third-party assessment under the relevant sectorial legislation

**and for AI systems with transparency risks**

**Implication:**

**medical devices that are or that contain software as safety component**
**and that are class IIa/B or higher are subject to AI Act**

# Legislative Lasagna

**AI SUMMIT**
COLUMBUS, OH • OCTOBER 25–27, 2022

insulin

- EU Data Act*
- EU Health Data Space*
- Artificial Intelligence Act*
- Radio Equipment Directive
- Machinery Directive
- Medical Device Regulation

Today, technical documentation of a closed-loop insulin pump needs to demonstrate compliance with three different legislations before the CE-mark can be assigned.

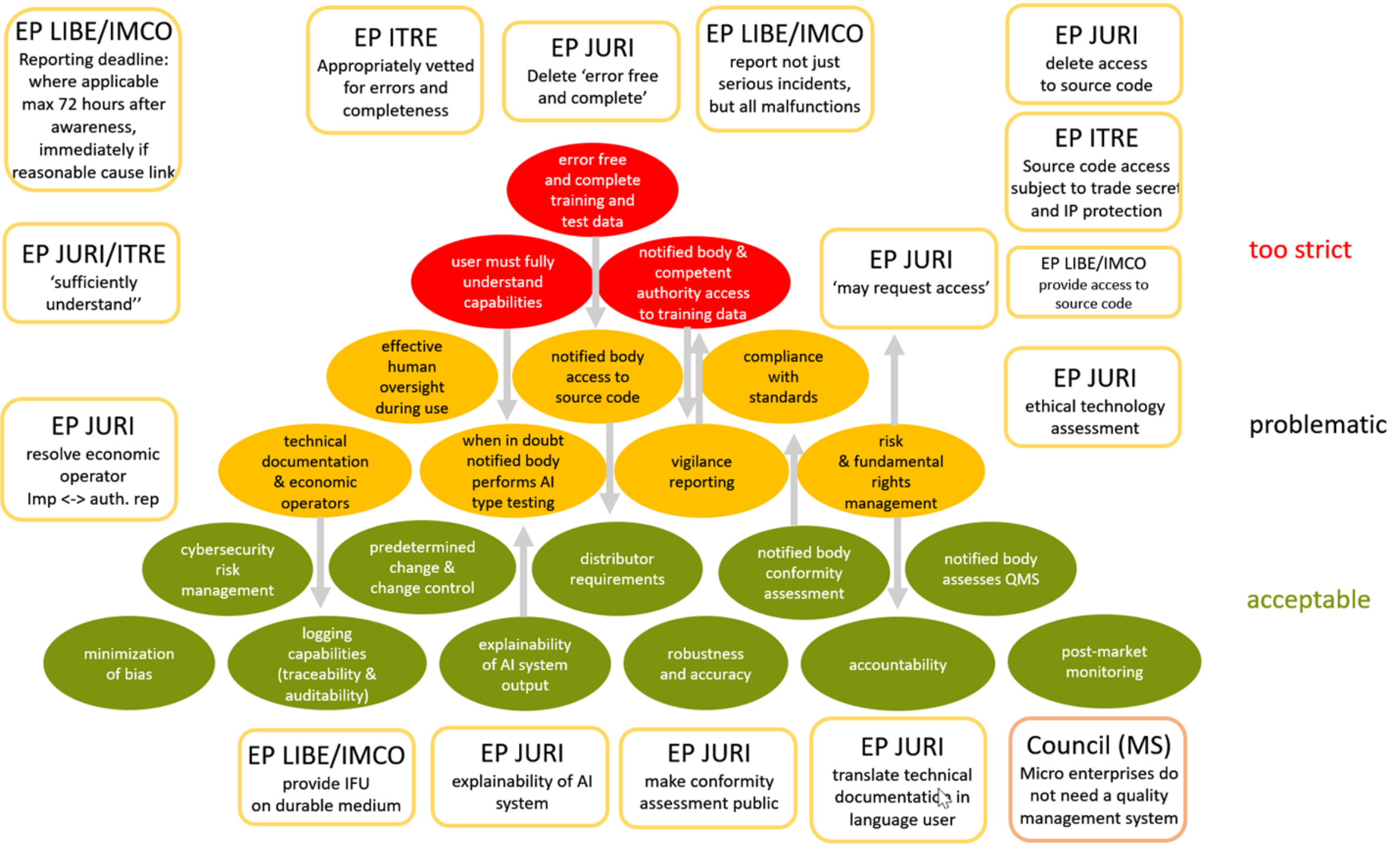*In the future, three additional legislations may come on top.

AFDO RAPS | HEALTHCARE PRODUCTS COLLABORATIVE | AFDO | RAPS

**too strict**

**problematic**

**acceptable**

EP LIBE/IMCO
Reporting deadline: where applicable max 72 hours after awareness, immediately if reasonable cause link

EP ITRE
Appropriately vetted for errors and completeness

EP JURI
Delete 'error free and complete'

EP LIBE/IMCO
report not just serious incidents, but all malfunctions

EP JURI
delete access to source code

EP ITRE
Source code access subject to trade secret and IP protection

EP JURI/ITRE
'sufficiently understand''

EP JURI
'may request access'

EP LIBE/IMCO
provide access to source code

EP JURI
resolve economic operator
Imp <-> auth. rep

EP JURI
ethical technology assessment

error free and complete training and test data

user must fully understand capabilities

notified body & competent authority access to training data

effective human oversight during use

notified body access to source code

compliance with standards

technical documentation & economic operators

when in doubt notified body performs AI type testing

vigilance reporting

risk & fundamental rights management

cybersecurity risk management

predetermined change & change control

distributor requirements

notified body conformity assessment

notified body assesses QMS

minimization of bias

logging capabilities (traceability & auditability)

explainability of AI system output

robustness and accuracy

accountability

post-market monitoring

EP LIBE/IMCO
provide IFU on durable medium

EP JURI
explainability of AI system

EP JURI
make conformity assessment public

EP JURI
translate technical documentation in language user

Council (MS)
Micro enterprises do not need a quality management system

# Standards with high operationalization value

for implementing AI Act requirements

**Overlaps & Conflicts:**     **extra costs, for little or no added value**

- **ISO/IEC 4213** Information technology — Artificial Intelligence — Assessment of ML classification performance
- **ISO/IEC 5259-3** Data quality for analytics and ML — Part 3: Data quality management requirements and guidelines

**IEC 62304 <->**
- **ISO/IEC 5338** Information technology — Artificial intelligence — AI system life cycle processes
- **ISO/IEC 5469** Artificial intelligence — Functional safety and AI systems

**ISO 14971 <->**
- **ISO/IEC 23894-2** Information Technology — Artificial Intelligence — Risk Management
- **ISO/IEC 24027** Information technology — Artificial intelligence (AI) — Bias in AI systems and AI aided decision making
- **ISO IEC 24029-1** Artificial Intelligence (AI) — Assessment of the robustness of neural networks — Part 1: Overview
- **ISO/IEC 38507** Information technology — Governance of IT — Governance implications of the use of artificial intelligence by organizations

**ISO 13485 <->**
- **ISO/IEC 42001** Information Technology — Artificial intelligence — Management system

List complied by AI Watch, joint initiative of European Commission and EC Joint Research Council
Above listed ISO/IEC SC42 standards are still under development

**AI SUMMIT**
COLUMBUS, OH • OCTOBER 25–27, 2022

# Lessons from Industry

## Where are we headed?

# Artificial Intelligence Principles @Google

1. Be socially beneficial.

2. Avoid creating or reinforcing unfair bias.

3. Be built and tested for safety.

4. Be accountable to people.

5. Incorporate privacy design principles.

6. Uphold high standards of scientific excellence.

7. Be made available for uses that accord with these principles.

ai.google/principles/

AFDO RAPS HEALTHCARE PRODUCTS COLLABORATIVE | AFDO RAPS

# Which includes things we will <u>not</u> do

We will not pursue certain AI applications…

likely to cause overall harm

weapons or those that direct injury

surveillance violating internationally accepted norms

purpose contravenes international law and human rights

ai.google/principles/

AFDO RAPS | HEALTHCARE PRODUCTS COLLABORATIVE | AFDO RAPS

Figure 1: Taxonomy of evaluation approaches for interpretability

https://arxiv.org/pdf/1702.08608.pdf

**Model Cards for Model Reporting**

https://arxiv.org/pdf/1810.03993.pdf



**Model Card**

- **Model Details.** Basic information about the model.
  - Person or organization developing model
  - Model date
  - Model version
  - Model type
  - Information about training algorithms, parameters, fairness constraints or other applied approaches, and features
  - Paper or other resource for more information
  - Citation details
  - License
  - Where to send questions or comments about the model
- **Intended Use.** Use cases that were envisioned during development.
  - Primary intended uses
  - Primary intended users
  - Out-of-scope use cases
- **Factors.** Factors could include demographic or phenotypic groups, environmental conditions, technical attributes, or others listed in Section 4.3.
  - Relevant factors
  - Evaluation factors
- **Metrics.** Metrics should be chosen to reflect potential real-world impacts of the model.
  - Model performance measures
  - Decision thresholds
  - Variation approaches
- **Evaluation Data.** Details on the dataset(s) used for the quantitative analyses in the card.
  - Datasets
  - Motivation
  - Preprocessing
- **Training Data.** May not be possible to provide in practice. When possible, this section should mirror Evaluation Data. If such detail is not possible, minimal allowable information should be provided here, such as details of the distribution over various factors in the training datasets.
- **Quantitative Analyses**
  - Unitary results
  - Intersectional results
- **Ethical Considerations**
- **Caveats and Recommendations**

Figure 1: Summary of model card sections and suggested prompts for each.

# Explainable AI

## Understand AI output and build trust

Explainable AI is a set of tools and frameworks to help you understand and interpret predictions made by your machine learning models, natively integrated with a number of Google's products and services. With it, you can debug and improve model performance, and help others understand your models' behavior. You can also generate feature attributions for model predictions in AutoML Tables, BigQuery ML and Vertex AI, and visually investigate model behavior using the What-If Tool.

AFDO RAPS | HEALTHCARE PRODUCTS COLLABORATIVE | AFDO RAPS

# What - If Tool

Interactive, visual debugging of black - box models

[Website](Website)

Probe classification and regression models, performing what - if analysis and analyzing fairness.

# Language Interpretability Tool

**Interactive, extensible, visual debugging of NLP models and beyond**

[Website](#)

Successor to the What - If Tool.

Probe models of all types (with a focus on NLP), explore model internals, prediction explanations, fairness, counterfactual generation, and more.

# Know Your Data

is an ML-based dataset exploration tools for rich, unstructured data

- Automatically computes signals

- Surface most biased data feature automatically through sorting and coloring

# Fairness Indicators

[Open-source library](#) that enables users to evaluate model performance for specific user groups ("sliced" analysis):

- Comes pre-loaded with common fairness metrics

- Provides interactive dashboard for rapid analysis & sharing insights with others

- Run analyses and visualize results in Jupyter notebooks or as part of TFX pipelines

## Fairness Indicators dashboard



**Select from common fairness metrics** one or multiple

**Set different decisions threshold(s)**

**Hover over metrics to inspect their definition**

Select metrics to display:
- ☐ Select all
- ☑ false_negative_rate
- ☐ false_positive_rate
- ☐ negative_rate
- ☐ positive_rate
- ☐ true_negative_rate
- ☐ true_positive_rate
- ☐ auc
- ☐ false_discovery_rate@0.9
- ☐ false_omission_rate@0.9

Baseline
race1:asian

fairness_indicators_metrics/false_negative_rate ⚙

Thresholds  0.9     Sort by  Slice

**Select a baseline** to compare performance against

**Choose slice(s)** to display (e.g. subgroups)

**View results** for your selections, side-by-side

# Model Remediation Library

[Open-source library](#) that enables users to train classifiers that equalize performance (provide "equal treatment") across a dimension, e.g. demographic group

Based on MinDiff modeling method (paper: *[Toward a better trade-off between performance and fairness with kernel-based distribution matching](#)* )

# Model Cards

[Model Cards](#) offer a transparency framework for organizing & communicating key information about a model in a standardized way.

Open-source [Model Card Toolkit](#) library facilitates and streamlines the creation of model cards

[Model Cards for Model Reporting](#) paper (2019)



## Model Card Toolkit

The Model Card Toolkit (MCT) library streamlines and automates generation of Model Cards, machine learning documents that provide context and transparency into a model's development and performance. Integrating the Model Card Toolkit into your ML pipeline will allow you to share your model's metadata and metrics with researchers, developers, reporters, and more.

MCT stores model card fields using a JSON schema. MCT can automatically populate those fields for TFX users via ML Metadata (MLMD). Model card fields can also be manually populated via a Python API. Some use cases of model cards include:

- Facilitating the exchange of information between model builders and product developers.
- Informing users of ML models to make better-informed decisions about how to use them (or how not to use them).
- Providing model information required for effective public oversight and accountability.

```
import model_card_toolkit

# Initialize the Model Card Toolkit with a path
model_card_output_path = ...
mct = model_card_toolkit.ModelCardToolkit(model

# Initialize the model_card_toolkit.ModelCard,
model_card = mct.scaffold_assets()
model_card.model_details.name = 'My Model'

# Write the model card data to a JSON file
mct.update_model_card_json(model_card)

# Return the model card document as an HTML pag
html = mct.export_format()
```

RAPS PRODUCTS COLLABORATIVE | AFDO RAPS

# Data Cards

Data Cards offer a structured way to document datasets & facilitate informed decision making for various stakeholders.

The [Data Cards Playbook](#) is a people-centered resource to help teams create customizable dataset documentation.

# DISCUSSION

**BUILDING TRUST IN AI SYSTEMS:** WHERE
ARE WE NOW?