

# Automated control and inspection of AI systems

*In accordance with regulatory frameworks: EU AI Act, ISO 42001, Greek Law 4961/2022*

Dr. Petros Stavroulakis

Trustworthy AI Consultant, code4thought

Visiting Professor of Computer Science, Evelpidon Military Academy

# About code4thought

## Origins

Founded in 2017 with a unique purpose, to render technology transparent, for large-scale software and AI-based systems.



## Location

Delivery center in Athens  
R&D center in Patras

## Software Quality

We evaluate and monitor the **non-functional** quality of any software system at every stage of their life cycle, through the analysis of source code and architecture with the SIGRID platform of our partner SIG



## Powered by





- \_decisively investing in **own R&D**
- \_international partnerships & project-participation
- \_strong ties with Academia
- \_ **TIER1 customer base**
- \_team with academic & enterprise experience
- \_participations and publications in international venues
- \_cross domain & technology agnostic solutions
- \_factual and well-researched approach

## Trustworthy AI

We help organizations govern any type of AI-based system at every stage of their life cycle, by **testing** & auditing AI models & their datasets with iQ4AI, our own platform



# Contents:

-  What is AI , really?
-  Regulatory Frameworks and Responsible AI
-  Compliance - centric AI measurement
-  Live Demo of iQ4AI model evaluation tool

# What is AI , really ?

*Basics of Artificial Intelligence and Machine Learning*

Dr. Petros Stavroulakis

Trustworthy AI Consultant, code4thought

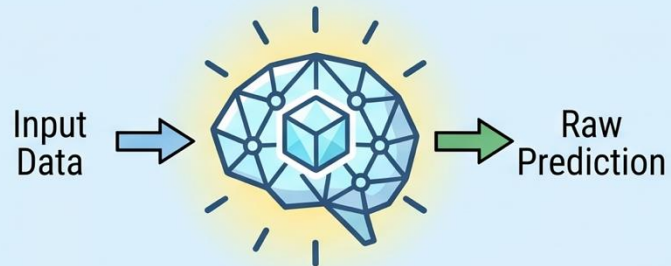
Visiting Professor of Computer Science, Evelpidon Military Academy

# AI models vs AI systems:

The 'brain'

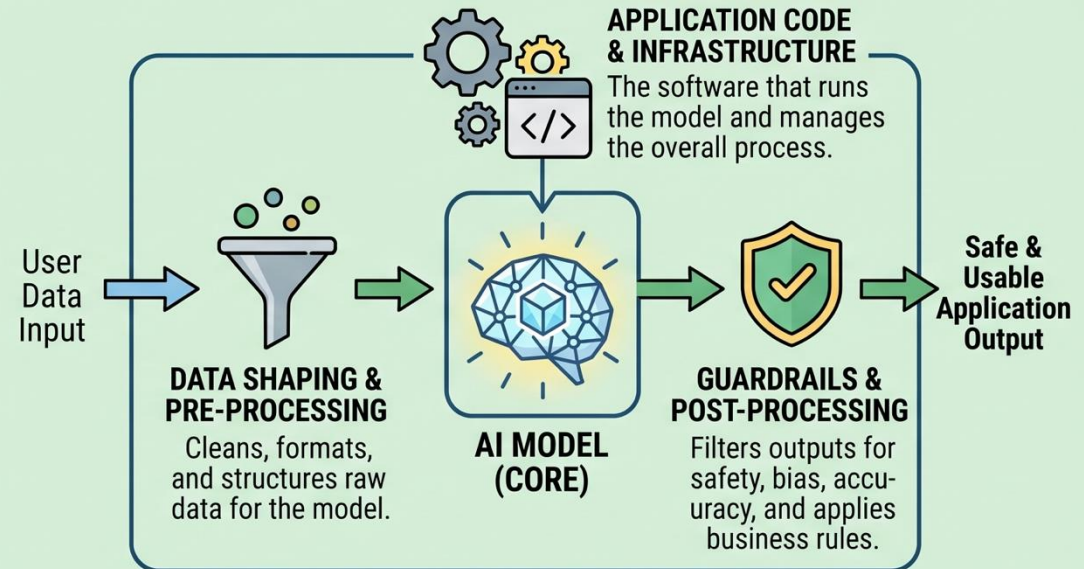
The 'body'

## AI MODEL: THE CORE



- **THE MATHEMATICAL ENGINE:** The core algorithm (e.g., weights, parameters).
- **MAKES PREDICTIONS:** Takes shaped input data and provides raw outputs based on trained patterns.
- **ISOLATED COMPONENT:** The fundamental reasoning part, not a full application.

## AI SYSTEM: THE COMPLETE SOLUTION



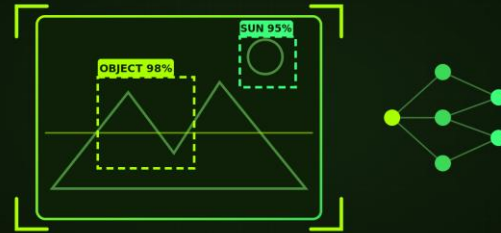
**THE MODEL IS THE REASONING ENGINE. THE SYSTEM IS THE COMPLETE WORKFLOW AND APPLICATION THAT MAKES IT USEFUL AND SAFE.**

# Examples of different types of AI models



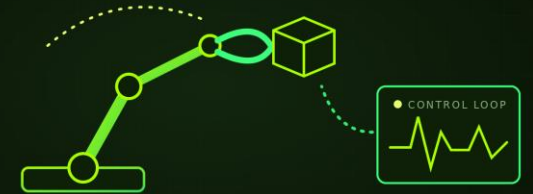
## SPEECH RECOGNITION

VOICE → TEXT



## IMAGE RECOGNITION

PIXELS → UNDERSTANDING



## MANIPULATION & CONTROL

DECISIONS → ACTIONS



## SOUND GENERATION

DATA → AUDIO



## IMAGE GENERATION

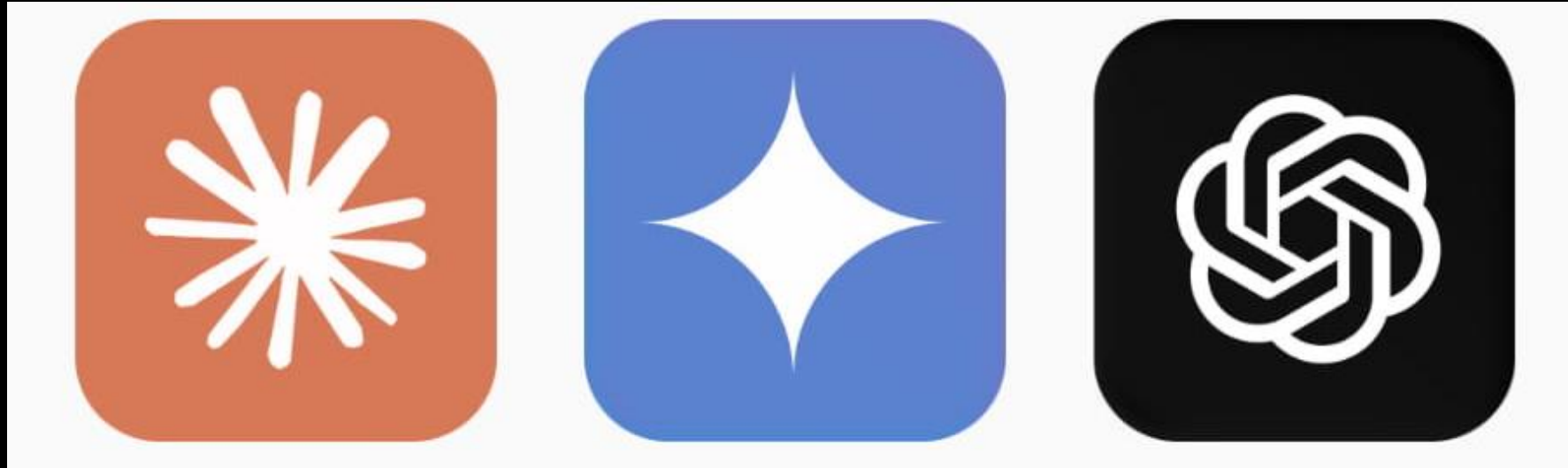
PROMPT → PICTURE



## TEXT GENERATION

PROMPT → WORDS

# Examples of different types of 'AI systems':



Claude

Gemini

ChatGPT

# code4thought - Automated document document compliance accelerator:

## Compliance Assessment Report

ISO Standard: ISO/IEC 42001-2023 (Sections scored: 7 / 20)

Date: 04 June 2026/

### Documents evaluated:

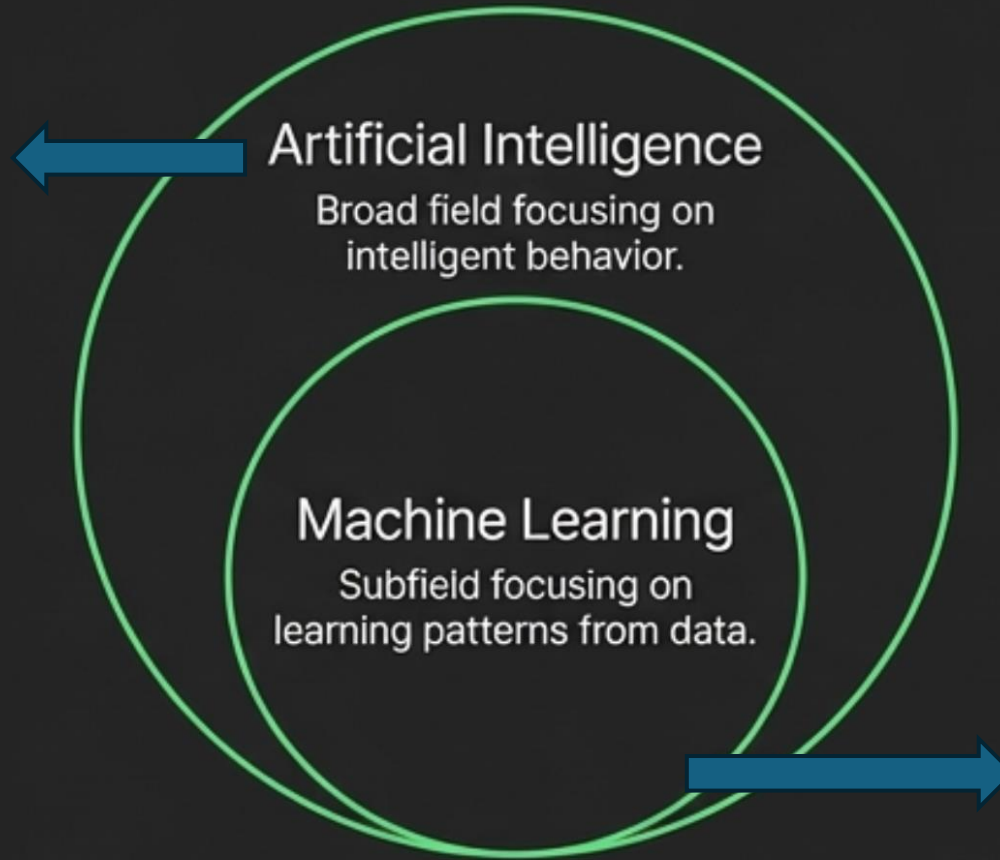
1. [Risk Matrix](#),
2. [Διαδικασία Αξιολόγησης & Έγκρισης Εργαλείων](#),
3. [Διαδικασία Αξιολόγησης Κινδύνων & Επιπτώσεων Χρήσης TN](#),
4. [Κατάλογος Αξιολόγησης Εργαλείων TN](#)

Average score: **3.0/5**



Section	Title	Score	What Exists	Gap	References
4.1	Understanding the organization and its context	2	Ο οργανισμός δηλώνει ρητά ότι δεν αναπτύσσει, δεν εκπαιδεύει και δεν διανέμει συστήματα τεχνητής νοημοσύνης σε τρίτους, αλλά χρησιμοποιεί αποκλειστικά έτοιμες υπηρεσίες TN τρίτων για εσωτερικούς λειτουργικούς σκοπούς, προσδιορίζοντας έτσι εν μέρει τον ρόλο του ως πελάτη/χρήστη/φορέα εφαρμογής TN.	Δεν εντοπίζεται σε κανένα έγγραφο επίσημος και δομημένος προσδιορισμός των εξωτερικών και εσωτερικών ζητημάτων, όπως απαιτείται από το ISO 42001· δεν υπάρχει ρητή ανάλυση πλαισίου τύπου SWOT, PESTLE ή ισοδύναμο τεκμήριο ανάλυσης.	Η γενική πολιτική διαχείρισης χρήσης TN ορίζει στο πεδίο εφαρμογής της ότι «ο οργανισμός δεν αναπτύσσει, δεν εκπαιδεύει και δεν διαθέτει συστήματα τεχνητής νοημοσύνης προς τρίτους και χρησιμοποιεί αποκλειστικά έτοιμες υπηρεσίες TN τρίτων μερών για εσωτερικούς επιχειρησιακούς σκοπούς».

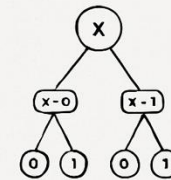
# Artificial intelligence vs Machine Learning



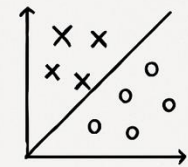
Learns patterns from data !

## EXPLORING MACHINE LEARNING MODELS

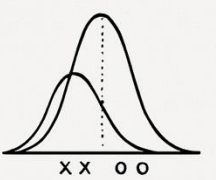
LOGICAL, GEOMETRIC, AND PROBABILISTIC APPROACHES



LOGICAL



GEOMETRIC



PROBABILISTIC

All ML is AI, but not all AI uses ML.

# Machine Learning vs Traditional Programming

## The Paradigm Shift

Traditional Programming:	[Inputs]	+	[Rules]	→	Output
Machine Learning:	[Inputs]	+	[Outputs]	→	Model (learned rules)

Minimizing the error between input and output

Main difference:

**Traditional Programming** – predictable engineered behavior

**Machine learning** - behavior depends on dataset and optimization procedure  
(requires continuous monitoring)

# Example

- Trying to "extract" rule  $y = \sin(x)$ , only by using a lookup table.

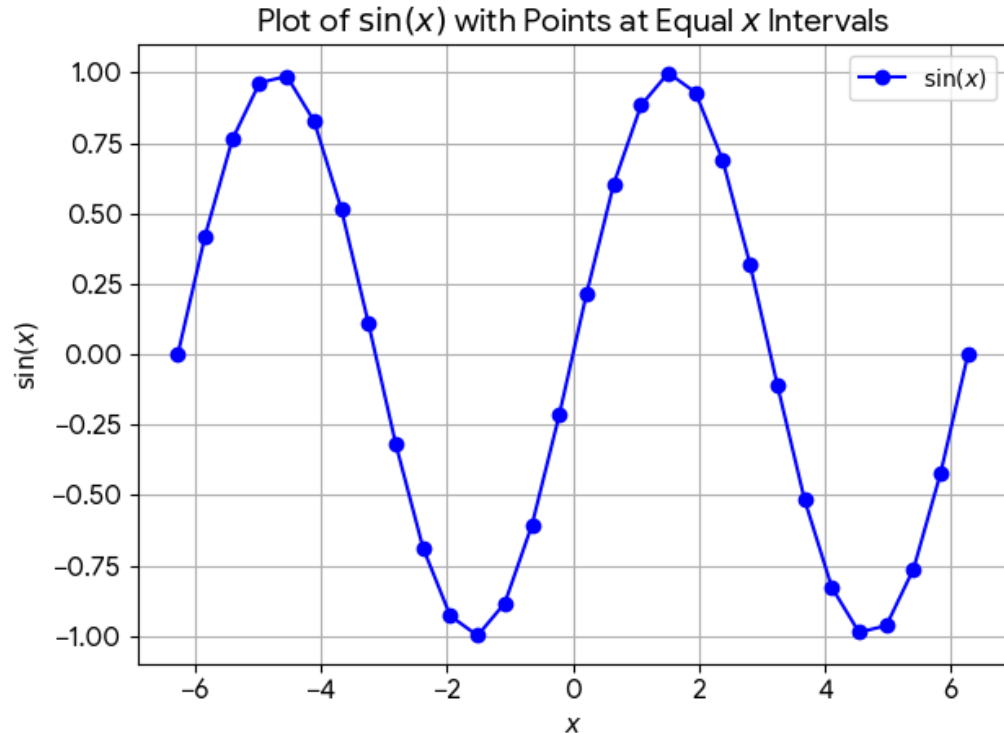


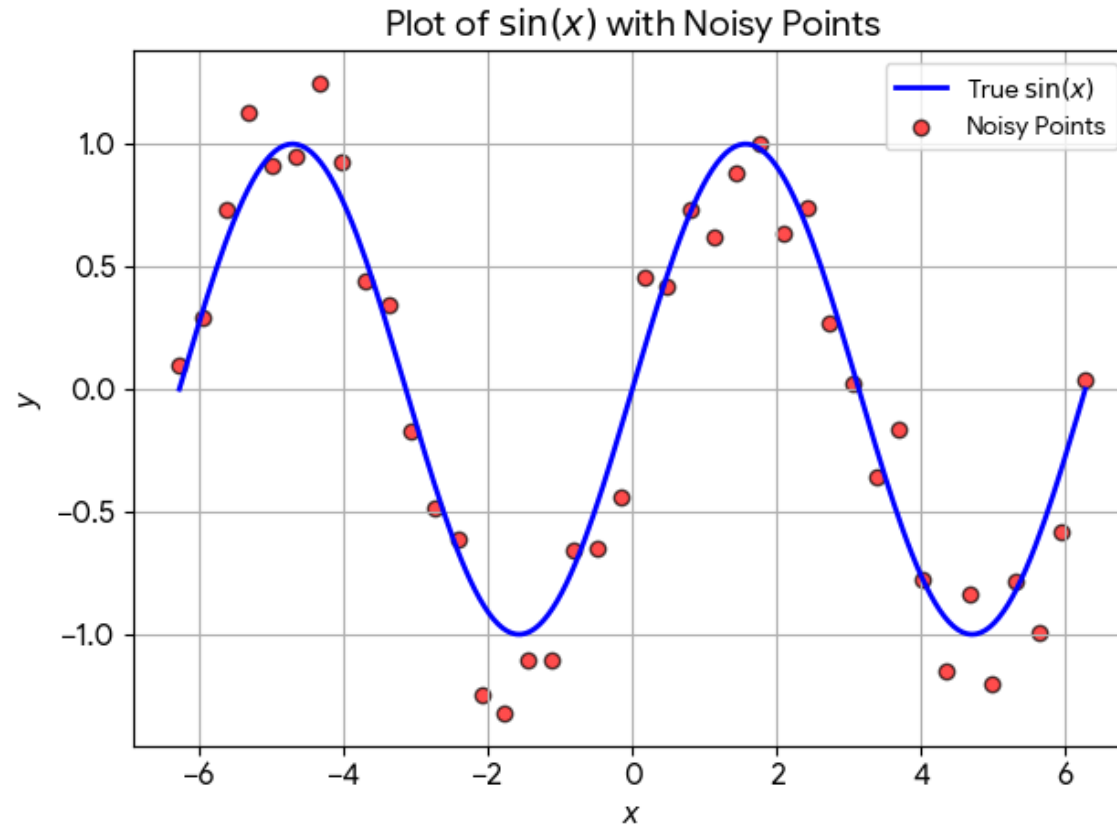
Table of  $\sin(\text{angle})$

Angle	sin (a)	Angle	sin (a)	Angle	sin (a)	Angle	sin (a)
0.0	0.0	25.0	.4226	46.0	.7193	71.0	.9455
1.0	.0174	26.0	.4384	47.0	.7314	72.0	.9511
2.0	.0349	27.0	.4540	48.0	.7431	73.0	.9563
3.0	.0523	28.0	.4695	49.0	.7547	74.0	.9613
4.0	.0698	29.0	.4848	50.0	.7660	75.0	.9659
5.0	.0872	30.0	.5000	51.0	.7772	76.0	.9703
6.0	.1045	31.0	.5150	52.0	.7880	77.0	.9744
7.0	.1219	32.0	.5299	53.0	.7986	78.0	.9781
8.0	.1392	33.0	.5446	54.0	.8090	79.0	.9816
9.0	.1564	34.0	.5592	55.0	.8191	80.0	.9848
10.0	.1736	35.0	.5736	56.0	.8290	81.0	.9877
11.0	.1908	36.0	.5878	57.0	.8387	82.0	.9903
12.0	.2079	37.0	.6018	58.0	.8480	83.0	.9926
13.0	.2249	38.0	.6157	59.0	.8571	84.0	.9945
14.0	.2419	39.0	.6293	60.0	.8660	85.0	.9962
15.0	.2588	40.0	.6428	61.0	.8746	86.0	.9976
16.0	.2756	41.0	.6561	62.0	.8829	87.0	.9986
17.0	.2924	42.0	.6691	63.0	.8910	88.0	.9994
18.0	.3090	43.0	.6820	64.0	.8988	89.0	.9998
19.0	.3256	44.0	.6947	65.0	.9063	90.0	1.00
20.0	.3420	45.0	.7071	66.0	.9135		
21.0	.3584			67.0	.9205		
22.0	.3746			68.0	.9272		
23.0	.3907			69.0	.9336		
24.0	.4067			70.0	.9397		

Use your browser "Print" command to make copies of this form.

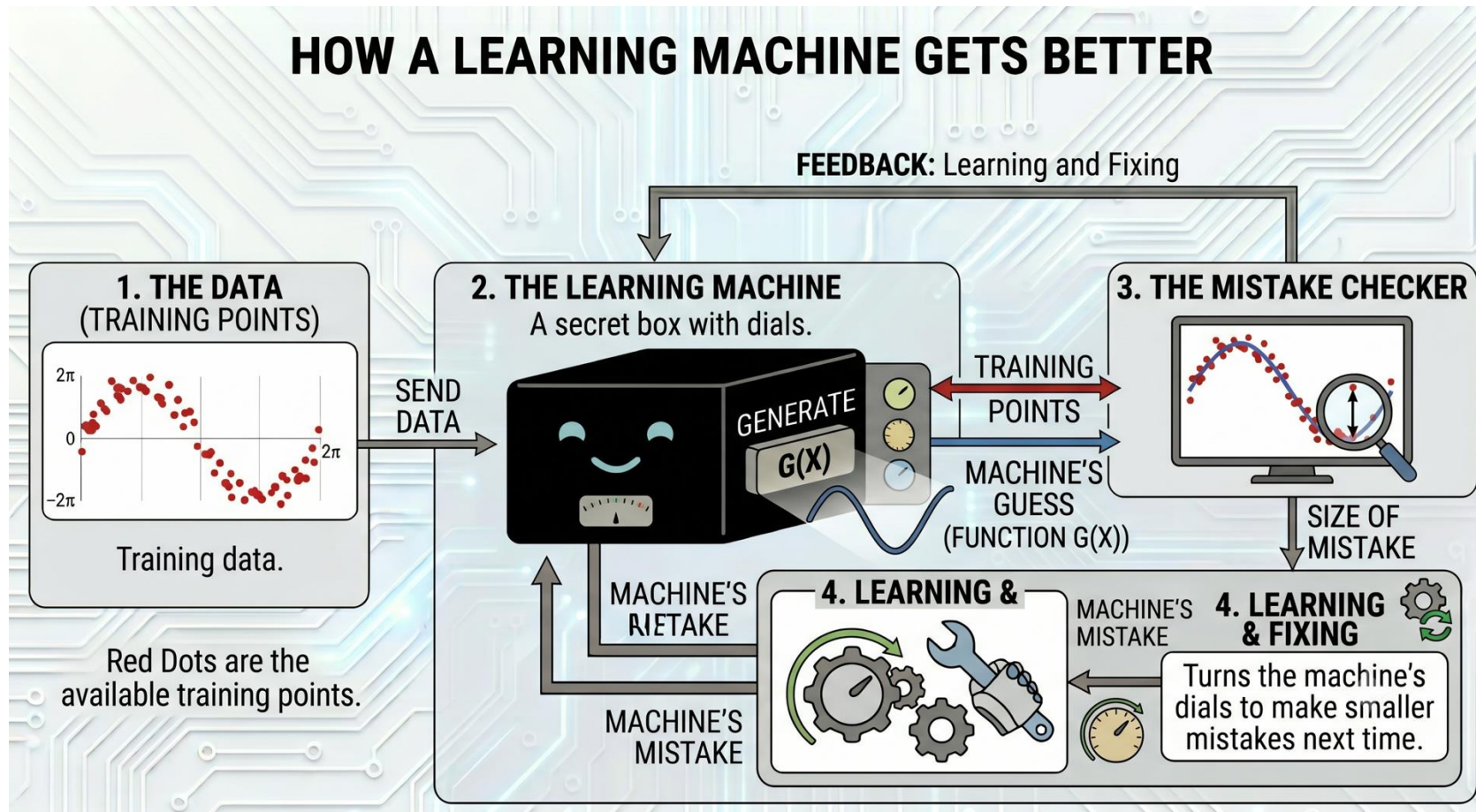
# Example

- In true life data contains 'noise' / poorly curated data:



# Example

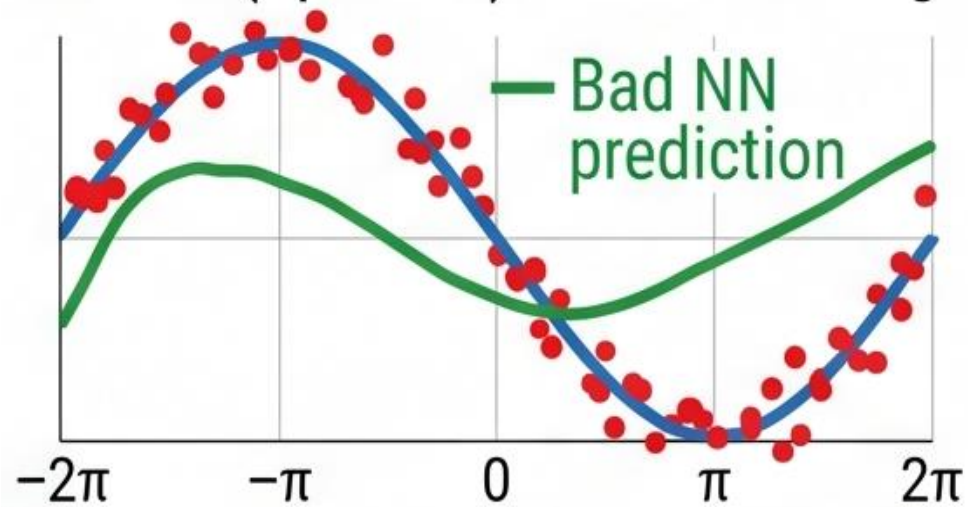
- Therefore the 'best we can do' is error minimization in multiple iterations:



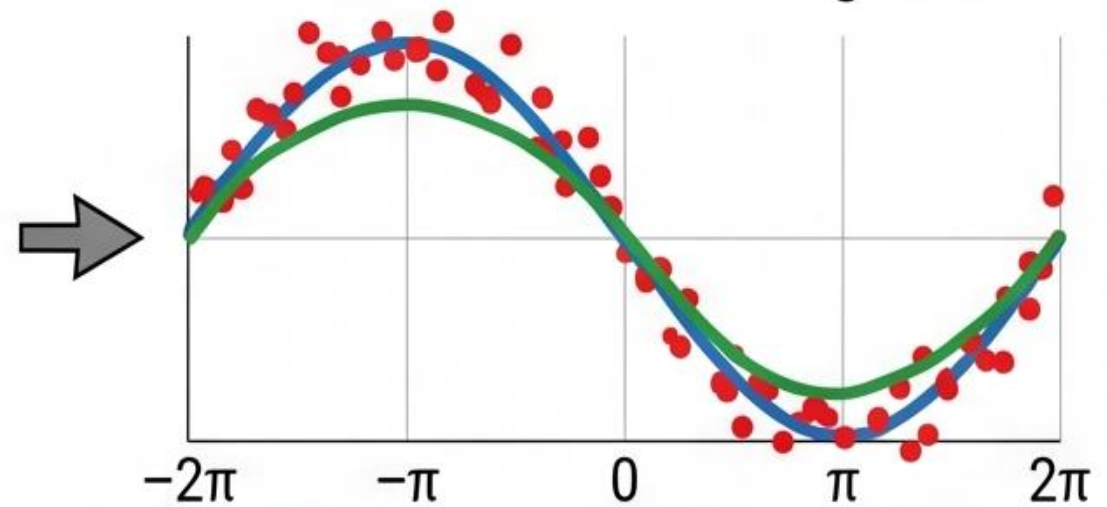
# Example

- Example of applying error minimization (green line) after some iterations to noisy input data (red), blue is the perfectly optimized result :

**START** (Epoch 0): Random Weights



**MIDDLE**: Minimizing Error



# The need for continuous monitoring:

Continuous monitoring of AI systems is required due to:

- Data Drift ( data input is very different to training data ) Performance Aspects
- Concept Drift ( core relationship changes e.g. from sinusoid to linear)
- Bias ( from training data / training procedure ) Regulation Aspects
- Fairness ( positive decision rates for different demographics )
- Explainability ( when we get a decision , can we explain what led to it ?)

# Regulatory Frameworks and Responsible AI

*Mapping AI Evaluation Requirements for the EU AI Act and ISO/IEC42001*

Dr. Petros Stavroulakis

Trustworthy AI Consultant, code4thought

Visiting Professor of Computer Science, Evelpidon Military Academy



# Disclaimer:

**code4thought** does not provide legal advice and has never represented itself as doing so, either explicitly or implicitly. All information is provided for educational purposes only and should not be relied upon as legal advice.

# RESPONSIBLE AI (The How)

### 1. FAIRNESS & BIAS

**HISTORICAL DATA BIAS ANALYSIS**      **BIAS MITIGATION FILTER**

OVER-REPRESENTED

UNDER-REPRESENTED

EQUITABLE TRAINING DATA

- Analyze Training Data for Hidden Biases
- Implement Advanced Bias Mitigation Techniques
- Monitor System Output for Non-Discrimination

Ensure AI systems do not discriminate against specific individuals or demographic groups.

### 2. TRANSPARENCY & EXPLAINABILITY

**MODEL SPECIFICATIONS DOC**      **DECISION PATH GENERATOR**

**DATA PROVENANCE TRACKER**      **WHY CREDIT WAS DENIED**

- Disclose AI Operations and Decision Logic
- Create Clear, Human-Understandable Explanations
- Provide Detailed Audit Trails and Documentation

Design AI models so their operations, data, and decisions can be understood.

### 3. ACCOUNTABILITY

**AI GOVERNANCE BOARD**

**DEVELOPMENT LEAD**      **DEPLOYMENT MANAGER**      **MONITORING AUDITOR**

**RESPONSIBILITY CHECKLIST**      **ASSIGNABILITY ACCOUNTABILITY**

- Establish Clear Human Oversight for System
- Define Individual Responsibility Across Lifecycle
- Create Communication Channels for User Grievances

Establish clear human oversight and a way to assign responsibility for system outcomes.

### 4. PRIVACY

**DATA DE-IDENTIFICATION**      **CONSENT GATEWAY**      **SECURE DATA VAULT**      **MANAGE USER CONSENT**

- Practice Data Minimization (only essential data)
- Prioritize User Consent and Subject Rights
- Use Robust Data Anonymization and De-identification

Protect user data and respect user data autonomy.

### 5. SECURITY

**DATA POISONING**      **PROMPT INJECTION**

**ROBUSTNESS TESTING**      **CYBER ATTACK**

- Robust Defense Against Cyber Threats
- Harden Models Against Adversarial Manipulation
- Secure the Entire End-to-End Lifecycle

Ensure technology is robust against cyber threats and adversarial manipulation.

### 6. SAFETY & RELIABILITY

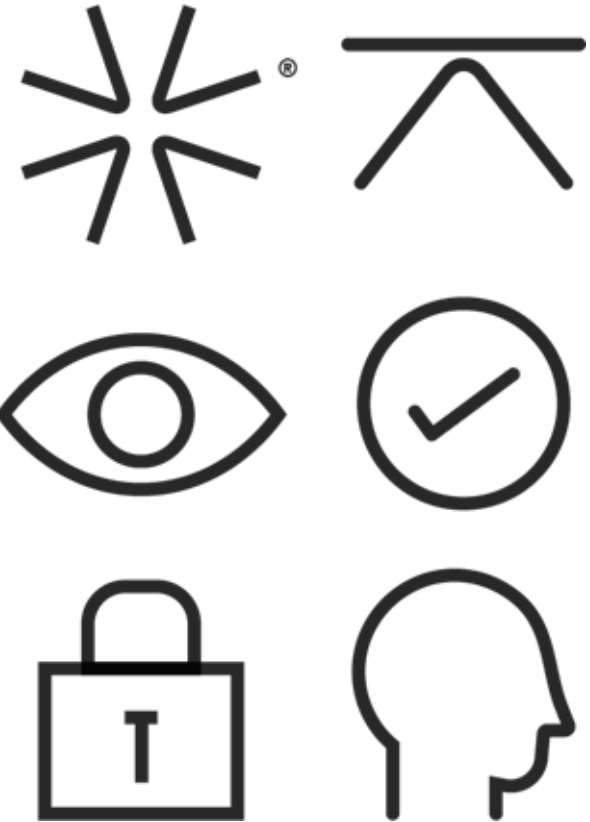
**Operator**

**POTENTIAL HARM**      **FINANCIAL CRASH**

**SAFE SHUTDOWN**      **SAFE-MODE ACTIVATION**      **PREVENTED HARM**

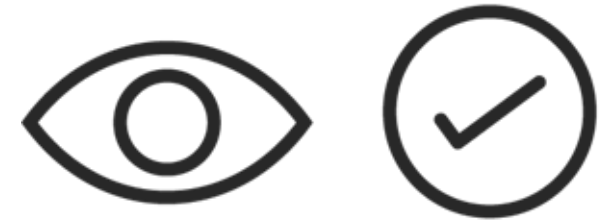
- Conduct Thorough Safety Risk Assessments
- Verify Model Accuracy and Generalization

Build accurate, dependable systems that do not cause harm.



**Definition:** The business governance and AI lifecycle processes by which produce Trustworthy AI

# TRUSTWORTHY AI (The result)



## 1. ACCOUNTABILITY

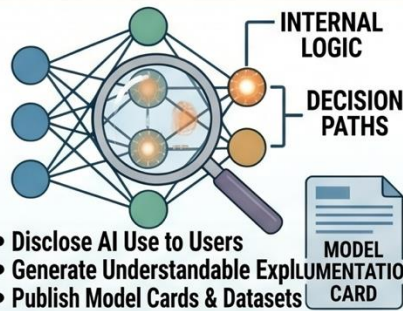


- Assign Clear System Ownership
- Maintain Detailed Audit Trails
- Establish Redress Mechanisms



Clear lines of responsibility for AI system actions and outcomes.

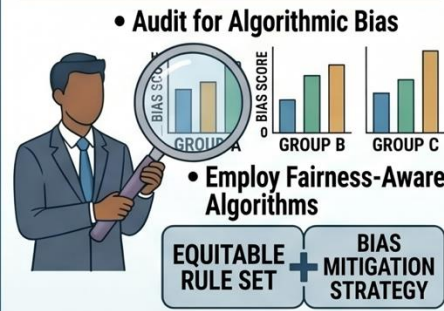
## 2. TRANSPARENCY



- Disclose AI Use to Users
- Generate Understandable Explanations
- Publish Model Cards & Datasets

Enable users to understand how and why AI decisions are made.

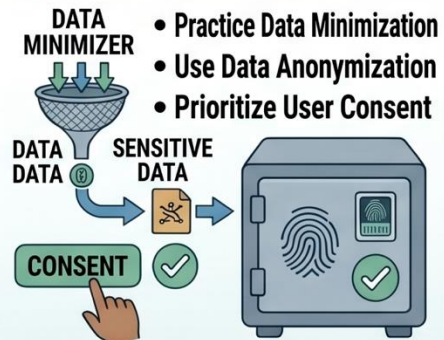
## 3. FAIRNESS



- Audit for Algorithmic Bias
- Employ Fairness-Aware Algorithms

Ensure equitable treatment and prevent discrimination across all user groups.

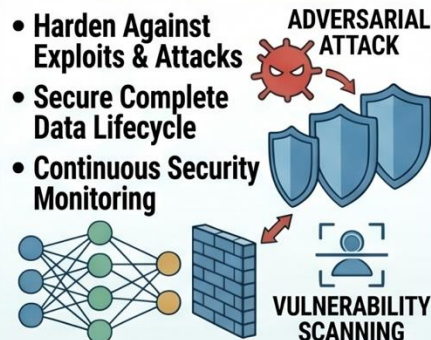
## 4. PRIVACY



- Practice Data Minimization
- Use Data Anonymization
- Prioritize User Consent

Protect sensitive personal data and respect user data autonomy.

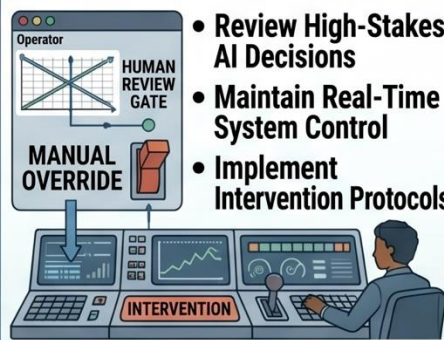
## 5. SECURITY



- Harden Against Exploits & Attacks
- Secure Complete Data Lifecycle
- Continuous Security Monitoring

Protect the AI system from external threats, manipulation, and vulnerabilities.

## 6. HUMAN OVERSIGHT

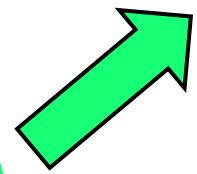


- Review High-Stakes AI Decisions
- Maintain Real-Time System Control
- Implement Intervention Protocols

Maintain effective human control and judgment over AI systems.

**Definition:** The result of Responsible AI, the produced system operates in a Trustworthy fashion.

# The vision:



**Disclaimer:** This image is used for educational and analysis purposes only under the Fair Use doctrine. The presentation does not own the copyright, nor does it imply an endorsement or promotion of the alcoholic beverage or brand shown."

**The Core Offering**  
The product beautifully and warmly lit, emphasizing craft, transparency, and care.

**The Mascot**  
A grounded, physical anchor that gives an inanimate product a living face and personality.

**The Hook**  
Bold, classic serif typography communicating uncompromising quality ('FULL of CHARACTER').



FULL of  
CHARACTER

Savour The Famous Grouse Responsibly.  
for the facts [drinkaware.co.uk](http://drinkaware.co.uk)

**The Mandate**  
The compliance text grounding the aspirational ad in ethical reality ('Savour... Responsibly').

[ Insert AI company logo. ]



FULL *of*  
INTEGRITY

Develop AI Responsibly.

Certified by:  code4thought®

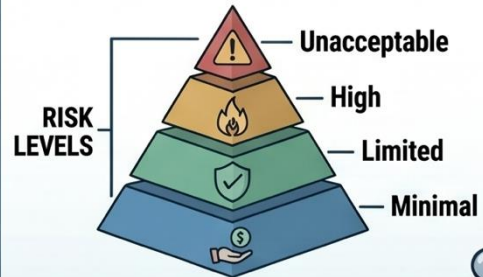


# EU AI Act

Regulation (EU) 2024/1689

# EU AI Act Intended Purpose

## 1. RISK-BASED REGULATION



Categorizes systems into four risk levels (Unacceptable, High, Limited, Minimal). Obligations are strictly proportionate to the potential level of harm to citizens.

## 2. FUNDAMENTAL RIGHTS PROTECTION



Ensures AI serves democracy and the rule of law. Specifically targets the mitigation of algorithmic bias and protects individual privacy and non-discrimination.

## 3. HUMAN OVERSIGHT & AGENCY



Mandates that humans remain "in the loop." High-risk AI must be designed for effective oversight by natural persons to prevent automated harms.

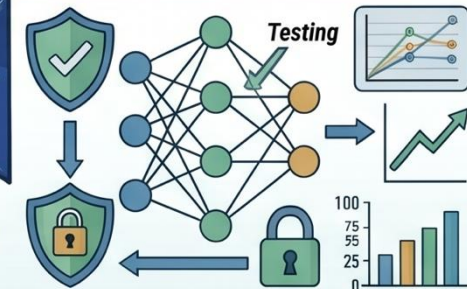
## 4. SINGLE MARKET HARMONIZATION



Replaces fragmented national laws with a unified legal framework across the EU to boost legal certainty and cross-border innovation for businesses.



## 5. TECHNICAL SAFETY & ROBUSTNESS



Strict requirements for accuracy, cybersecurity, and data quality. Model requirements for accuracy, cybersecurity, and data quality. Models must be resilient against adversarial attacks and technically dependable.

**The EU AI Act purpose:**

*Establishing the world's first comprehensive legal framework to ensure AI systems are , Safe , Trustworthy and Respectful of fundamental European Values*

# EU AI Act — Regulation (EU) 2024/1689

\_ **Entered into force August 2024;** phased application through 2027. Risk-based regulation covering AI placed on the EU market or affecting persons in the EU.

\_ **Evaluation requirements concentrate on high-risk AI systems** (Title III, Chapter 2) and on general-purpose AI models (Chapter V).

\_ **Six trustworthiness properties are explicitly required:** accuracy, robustness, cybersecurity, fairness/non-discrimination, transparency, explainability/interpretability.

\_ **Compliance is demonstrated through technical documentation** (Art. 11 + Annex IV), quality management (Art. 17), conformity assessment, and post-market monitoring (Art. 72).



# EU AI Act — Accuracy

**Article 15(1)** — high-risk AI systems shall be designed and developed to achieve an appropriate level of accuracy.

**Article 15(2)** — the Commission shall encourage the development of benchmarks and measurement methodologies.

**Article 15(3)** — accuracy levels and relevant accuracy metrics shall be declared in the accompanying instructions for use.

**Article 13(3)(b)(ii)** — transparency to deployers: instructions must include the metrics used to measure accuracy.

**Annex IV (2)(g)** — technical documentation must describe the metrics used to assess accuracy.

**Recital 74** — appropriate accuracy throughout the lifecycle; performance must be measurable and reproducible.



# EU AI Act — Robustness

**Article 15(1)** — high-risk AI systems shall achieve an appropriate level of robustness.

**Article 15(4) first subparagraph** — resilience against errors, faults and inconsistencies; **technical and organisational measures including redundancy solutions and fail-safe plans.**

**Article 15(4) second subparagraph** — addressing feedback loops where outputs influence future inputs, including bias amplification risks.

**Annex IV (2)(g)** — robustness metrics in the technical documentation, including stress-testing and behaviour under degraded conditions.

**Recital 75** — design for resilience, with emphasis on operating environment, foreseeable misuse and unexpected situations.



# EU AI Act — Cybersecurity

**Article 15(1) and 15(5)** — appropriate level of cybersecurity; resilience to attempts by unauthorised third parties to alter use, outputs or performance.

**Article 15(5) second subparagraph** — technical solutions for AI-specific vulnerabilities, including measures to prevent, detect, respond to, resolve and control: data poisoning, model poisoning, adversarial examples, model evasion, and confidentiality attacks.

**Annex IV (2)(h)** — cybersecurity measures put in place must be documented in the technical file.

**Article 72** — post-market monitoring shall include cybersecurity-relevant incidents and trigger updates where appropriate.

**Recitals 76–77** — cybersecurity obligations are required to be commensurate with risk and to integrate with the NIS2 / Cyber Resilience Act framework.



# EU AI Act — Fairness & Non-discrimination

**Article 10(2)(f)** — examination of training, validation and testing data in view of possible biases likely to affect health, safety or fundamental rights or to lead to discrimination.

**Article 10(2)(g)** — appropriate measures to detect, prevent and mitigate identified biases.

**Article 10(3)** — datasets must be relevant, sufficiently representative, free of errors and complete for the intended purpose.

**Article 10(5)** — permits processing of special categories of personal data strictly for bias detection and correction, subject to safeguards.

**Article 27** — Fundamental Rights Impact Assessment (FRIA) for deployers of certain high-risk AI systems.

**Recitals 67 and 70** — interpretive guidance on data quality, representativeness and discrimination.



# EU AI Act — Transparency

**Article 13** — transparency and provision of information to deployers: intended purpose, capabilities, limitations, expected accuracy and robustness, foreseeable misuse, performance on specified groups.

**Article 50** — transparency obligations for providers and deployers of certain AI systems: chatbots disclosing AI interaction, marking of synthetic / deepfake content, notification for emotion recognition and biometric categorisation.

**Article 53(1)(a)–(d)** — general-purpose AI model providers: technical documentation, information to downstream providers, and a publicly available summary of training data.

**Recitals 27, 72, 132–134** — principles of transparency and obligations regarding synthetic content.



# EU AI Act — Explainability & Interpretability

**Article 13(1)** — high-risk AI systems shall be designed and developed in such a way as to be sufficiently transparent to enable deployers to interpret a system's output and use it appropriately.

**Article 14(4)(c)** — natural persons assigned to human oversight must be able to correctly interpret the high-risk AI system's output.

**Article 86** — right to explanation of individual decision-making: affected persons subject to a decision based on a high-risk AI system may obtain clear and meaningful explanations of the role of the AI system in the decision-making procedure and the main elements of the decision taken.

**Recital 27** — the trustworthy-AI principle of explicability / intelligibility.

**Recital 171** — explanations to affected persons must be meaningful and proportionate.



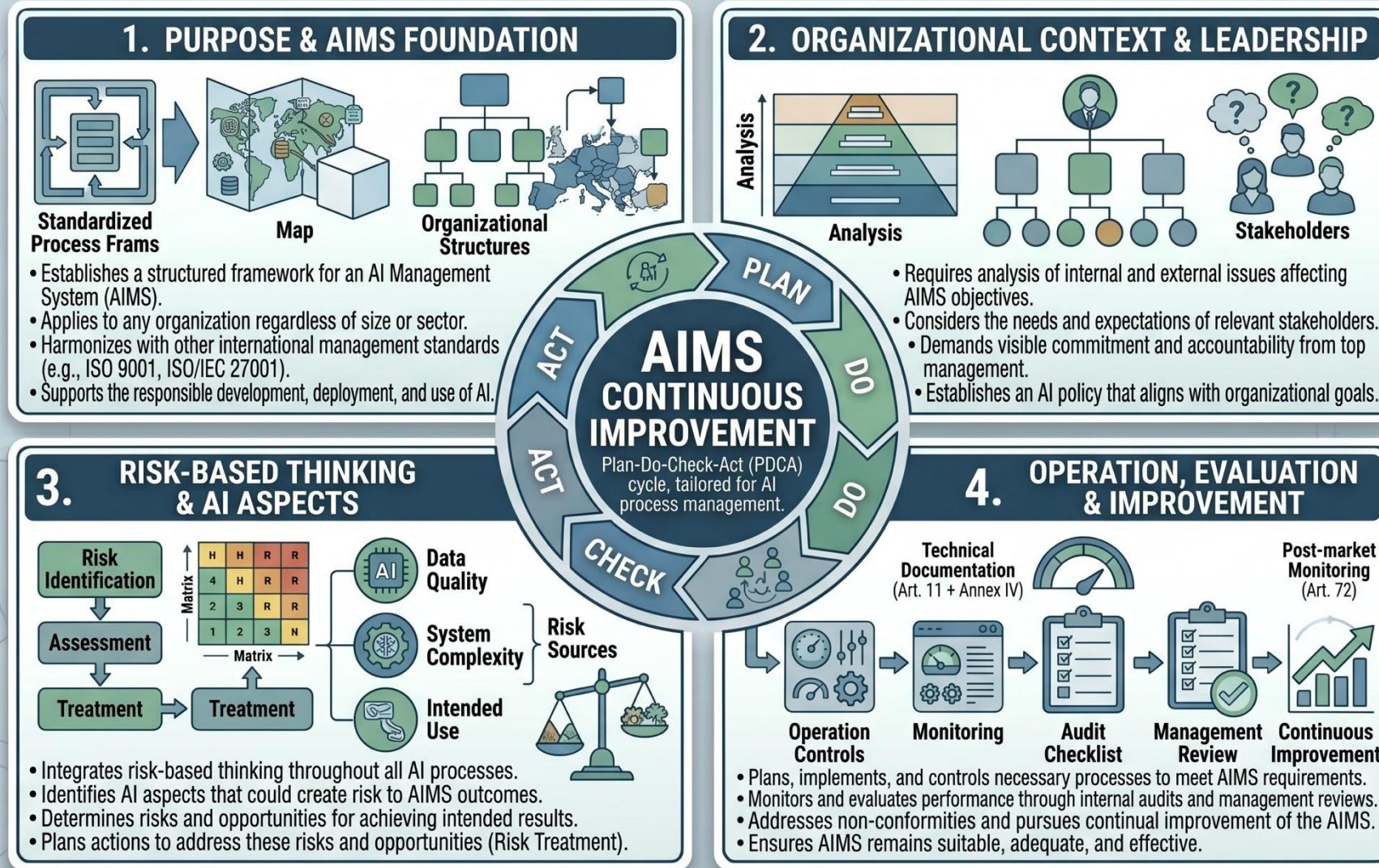


# ISO/IEC 42001:2023

AI Management System — clauses and Annex A controls

# ISO/IEC 42001:2023 — Main Purpose

Core principles of the Artificial Intelligence Management System (AIMS) standard, enabling organizations to manage AI-related processes effectively and responsibly.



**Core Purpose:** The establishment of an AI management system AIMS

# ISO/IEC 42001:2023 — Clauses related to AI evaluation

**Clause 6.1.2** — AI **risk assessment process**.

**Clause 6.1.3** — AI **risk treatment process**.

**Clause 6.1.4** — AI system **impact assessment** (covering individuals, groups, society).

**Clauses 8.2 / 8.3 / 8.4** — **operational implementation of risk assessment, risk treatment and impact assessment**.

**Clause 9.1** — **monitoring, measurement, analysis and evaluation of AIMS performance**.

**Clause 9.2** — **internal audit**; Clause 9.3 — management review.

**Clause 10** — **continual improvement and nonconformity management**.



# ISO/IEC 42001 — Annex A controls relevant to evaluation

**Impact assessment** — A.5.2 (process), A.5.3 (documentation), A.5.4 (impact on individuals/groups, fairness), A.5.5 (societal impacts).

**Lifecycle & testing** — A.6.1.2 objectives; A.6.1.3 processes; A.6.2.2 requirements & specification (accuracy, robustness, security, fairness, transparency, explainability); A.6.2.3 design documentation; A.6.2.4 verification & validation; A.6.2.5 deployment; A.6.2.6 operation & monitoring; A.6.2.7 technical documentation; A.6.2.8 event logs.

**Data quality** — A.7.2 data for development; A.7.3 acquisition; A.7.4 quality of data; A.7.5 provenance; A.7.6 preparation.

**Information to stakeholders** — A.8.2 documentation for users; A.8.3 external reporting; A.8.4 incident communication; A.8.5 information for interested parties.

**Responsible use** — A.9.2 processes; A.9.3 objectives; A.9.4 intended use.

**Third parties** — A.10.2 allocating responsibilities; A.10.3 suppliers; A.10.4 customers.





# Greek Law 4961/2022

«Αναδυόμενες τεχνολογίες πληροφορικής και επικοινωνιών, ενίσχυση της ψηφιακής διακυβέρνησης και άλλες διατάξεις.»

# Greek Law 4961/2022

Not the Greek mirror-law for EU AI Act , as has been suggested that has not yet been voted in.

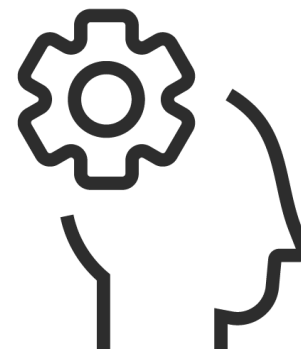
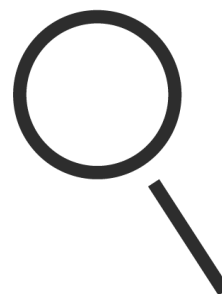
Greek Law 4961/2022 focuses on certain similar aspects, but was created before the EU AI Act:

- Defines how public and private sector uses and engages with AI systems
- Precursor to full EU AI Act (EU AI Act was voted in the EU in 2024)
- Algorithmic Impact Assessment (**Αλγοριθμική Εκτίμηση Αντικτύπου (ΑΙΑ)**) similar to EU AI Act's FRIA
- AI Systems Register (**Μητρώα Τεχνητής Νοημοσύνης για συστήματα υψηλού κινδύνου** ) similar to EU AI Act's European Registry
- Transparency and Accountability ( **Αρχές Διαφάνειας και Λογοδοσίας** )
- Bias (**Τήρηση της αρχής της ίσης μεταχείρισης και της καταπολέμησης των διακρίσεων**)
- Does not contradict GDPR

# Greek Law 4961/2022 – Safe and responsible use

## Αρθρο 1:

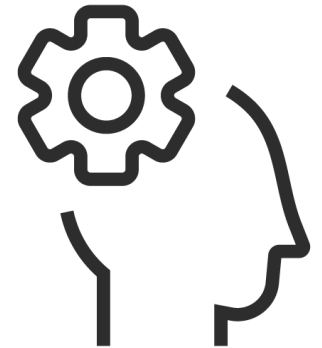
«Σκοπός του παρόντος Μέρους είναι η δημιουργία του κατάλληλου θεσμικού υποβάθρου για τη **θεμιτή και ασφαλή αξιοποίηση των δυνατοτήτων της τεχνολογίας τεχνητής νοημοσύνης από φορείς του δημόσιου και του ιδιωτικού τομέα** και η ενίσχυση της ανθεκτικότητας της δημόσιας διοίκησης απέναντι σε απειλές στον κυβερνοχώρο.»



# Greek Law 4961/2022 - Accountability

## Αρθρο 2:

«Αντικείμενο του παρόντος Μέρους είναι η θέσπιση ρυθμίσεων για τη διαμόρφωση των κατάλληλων **εγγυήσεων** για τη **διασφάλιση των δικαιωμάτων των φυσικών και των νομικών προσώπων** και την **ενίσχυση της λογοδοσίας και της διαφάνειας** κατά τη χρήση συστημάτων τεχνητής νοημοσύνης»

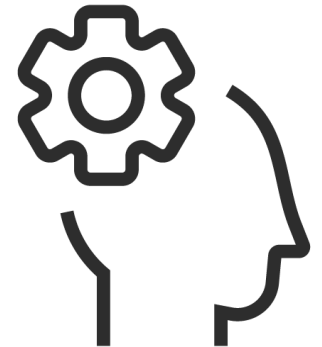
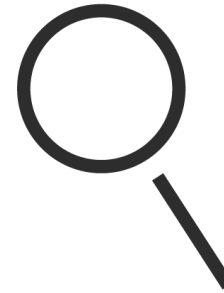


# Greek Law 4961/2022. – AI used in cases where citizen right protection guarantees have been provided

## Αρθρο 4:

1. Οι φορείς του δημόσιου τομέα κατά την έννοια της περ. α' της παρ. 1 του άρθρου 14 του ν. 4270/2014 (Α' 143) δύνανται, κατά την άσκηση των αρμοδιοτήτων τους, να χρησιμοποιούν συστήματα τεχνητής νοημοσύνης για τη διαδικασία λήψης ή την υποστήριξη της διαδικασίας λήψης μιας απόφασης ή την έκδοση πράξης, οι οποίες επηρεάζουν τα δικαιώματα ενός φυσικού ή νομικού προσώπου, **μόνον εφόσον η χρήση αυτή προβλέπεται ρητά σε ειδική διάταξη νόμου που περιλαμβάνει κατάλληλες εγγυήσεις για την προστασία των δικαιωμάτων αυτών.**

2. Από το πεδίο εφαρμογής του παρόντος κεφαλαίου **εξαιρούνται** τα Υπουργεία Εθνικής Άμυνας και Προστασίας του Πολίτη, οι εποπτευόμενοι φορείς αυτών, καθώς και η Εθνική Υπηρεσία Πληροφοριών (Ε.Υ.Π.).



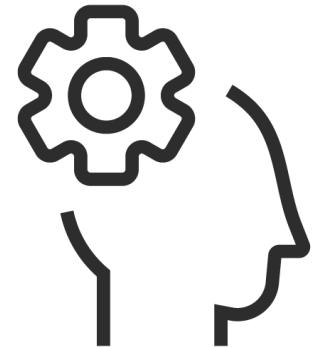
# Greek Law 4961/2022 – Impact assessment

## Άρθρο 5: Αλγοριθμική εκτίμηση αντικτύπου

1. Κάθε φορέας του δημόσιου τομέα που χρησιμοποιεί σύστημα τεχνητής νοημοσύνης της παρ. 1 του άρθρου 4, **πριν από την έναρξη λειτουργίας του συστήματος, εκπονεί αλγοριθμική εκτίμηση αντικτύπου.**
2. Κατά την εκπόνηση της αλγοριθμικής εκτίμησης αντικτύπου λαμβάνονται υπόψη, ιδίως, οι ακόλουθες πληροφορίες:

.... β) **οι δυνατότητες, τα τεχνικά χαρακτηριστικά και οι παράμετροι λειτουργίας του συστήματος,**

.... ε) **οι κίνδυνοι που ενδέχεται να προκύψουν για τα δικαιώματα, τις ελευθερίες και τα έννομα συμφέροντα των φυσικών ή νομικών προσώπων, στα οποία αφορά ή τα οποία επηρεάζει η λήψη της απόφασης** και ...



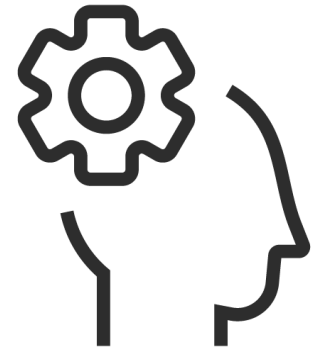
# Greek Law 4961/2022 - Transparency and Explainability

## Αρθρο 6 – Υποχρεώσεις διαφάνειας

« ... παρέχει, δημόσια, πληροφορίες σχετικά με:

- α) τον χρόνο έναρξης λειτουργίας του συστήματος,
- β) τις παραμέτρους λειτουργίας, τις δυνατότητες και τα τεχνικά χαρακτηριστικά του συστήματος,
- γ) τις κατηγορίες των αποφάσεων που λαμβάνονται ή των πράξεων που εκδίδονται με τη συμμετοχή του συστήματος ή υποστηρίζονται από αυτό και
- δ) τη διενέργεια αλγοριθμικής εκτίμησης αντικτύπου.»

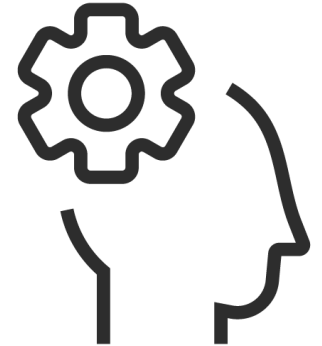
<< ότι το φυσικό ή νομικό πρόσωπο στο οποίο αφορά η λήψη της απόφασης ή η έκδοση της πράξης, λαμβάνει γνώση των παραμέτρων στις οποίες στηρίχθηκε η λήψη της απόφασης ή η έκδοση της πράξης σε κατανοητή και εύκολα προσβάσιμη μορφή, συμπεριλαμβανομένων των μορφών που διευκολύνουν την ενημέρωση των ατόμων με αναπηρία.>>



# Greek Law 4961/2022 - Transparency and Explainability

## Άρθρο 7 – Υποχρεώσεις αναδόχων συστημάτων τεχνητής νοημοσύνης

<<υποχρέωση του αναδόχου να παρέχει προς τον φορέα του δημόσιου τομέα τις πληροφορίες των περ. β' και γ' της παρ. 1 και της παρ. 2 του άρθρου 6, συμπεριλαμβανομένης της υποχρέωσης παραίτησής του από την άσκηση αξιώσεων που ενδέχεται να θέσουν σε κίνδυνο το δικαίωμα των φυσικών ή νομικών προσώπων για παροχή πληροφοριών, σύμφωνα με την παρ. 1 του άρθρου 6.>>



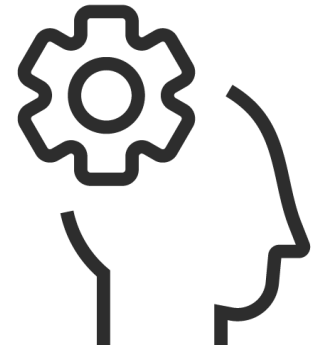
# Greek Law 4961/2022 - Safety

## Άρθρο 8

<<1. Κάθε φορέας του δημόσιου τομέα **υποχρεούται να τηρεί μητρώο με τα συστήματα τεχνητής νοημοσύνης της παρ. 1 του άρθρου 4 που χρησιμοποιεί**, το οποίο επικαιροποιεί μέχρι την 1η Μαρτίου κάθε έτους, και, σε κάθε περίπτωση, όταν θέτει σε λειτουργία νέο σύστημα.

...

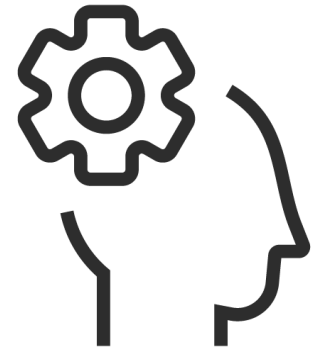
β) τον χρόνο έναρξης λειτουργίας του συστήματος, γ) τις παραμέτρους λειτουργίας, τις δυνατότητες και τα τεχνικά χαρακτηριστικά του συστήματος, δ) τις βασικές πληροφορίες του συστήματος, όπως η ονομασία, η έκδοση και η επωνυμία του κατασκευαστή, ε) **μέτρα που λαμβάνονται για την ασφαλή λειτουργία του,** >>



# Greek Law 4961/2022 - Bias

## Άρθρο 9

1. Κάθε επιχείρηση του ιδιωτικού τομέα, εφόσον χρησιμοποιεί σύστημα τεχνητής νοημοσύνης, το οποίο επηρεάζει οποιαδήποτε διαδικασία λήψης αποφάσεων σχετικά με τους εργαζομένους ή τους υποψήφιους εργαζομένους και έχει **αντίκτυπο στις συνθήκες εργασίας, την επιλογή, την πρόσληψη ή την αξιολόγησή τους**, σε κάθε περίπτωση πριν την πρώτη χρήση του, παρέχει επαρκή και σαφή πληροφόρηση σε κάθε εργαζόμενο ή υποψήφιο εργαζόμενο, η οποία περιλαμβάνει κατ' ελάχιστον τις παραμέτρους στις οποίες στηρίζεται η λήψη της απόφασης, με την επιφύλαξη των περιπτώσεων που προϋποθέτουν προηγούμενη ενημέρωση και διαβούλευση και **διασφαλίζει την τήρηση της αρχής της ίσης μεταχείρισης και της καταπολέμησης των διακρίσεων** στην απασχόληση και την εργασία λόγω φύλου, φυλής, χρώματος, εθνικής ή εθνοτικής καταγωγής γενεαλογικών καταβολών, θρησκευτικών ή άλλων πεποιθήσεων, αναπηρίας ή χρόνιας πάθησης, ηλικίας, οικογενειακής ή κοινωνικής κατάστασης, σεξουαλικού προσανατολισμού, ταυτότητας ή χαρακτηριστικών φύλου.

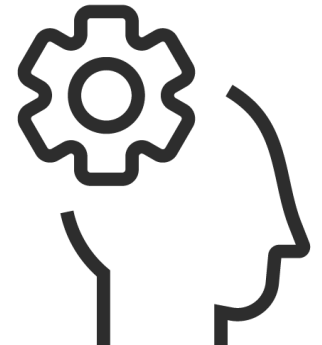


# Greek Law 4961/2022 – Human Profiling

## Άρθρο 10

1. Κάθε επιχείρηση του ιδιωτικού τομέα, η οποία συνιστά μεσαία ή μεγάλη οντότητα κατά την έννοια των παρ. 5 και 6, αντίστοιχα, του άρθρου 2 του ν. 4308/2014 (Α' 251) τηρεί, σε ηλεκτρονική μορφή, μητρώο των συστημάτων τεχνητής νοημοσύνης τα οποία χρησιμοποιεί είτε στο πλαίσιο κατάρτισης προφίλ καταναλωτών είτε στο πλαίσιο αξιολόγησης των πάσης φύσεως εργαζομένων της ή συνεργαζόμενων με αυτήν φυσικών προσώπων.

Στο μητρώο περιλαμβάνονται για κάθε σύστημα, κατ' ελάχιστον, οι εξής πληροφορίες: α) η περιγραφή των παραμέτρων λειτουργίας, των δυνατοτήτων και των τεχνικών χαρακτηριστικών του συστήματος,



# Compliance – centric AI measurement

*Automated measurement methods for AI systems*

Dr. Petros Stavroulakis

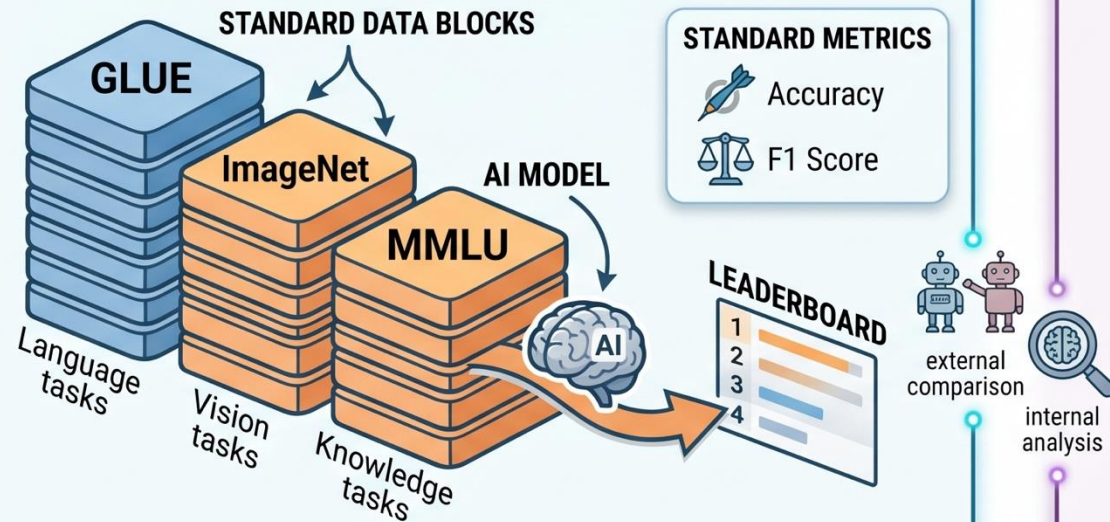
Trustworthy AI Consultant, code4thought

Visiting Professor of Computer Science, Evelpidon Military Academy

# Types of evaluation methods

## TWO DIFFERENT TYPES OF EVALUATION METHODS FOR TESTING AI MODELS

### 1. BENCHMARK DATASETS (Automated Testing)



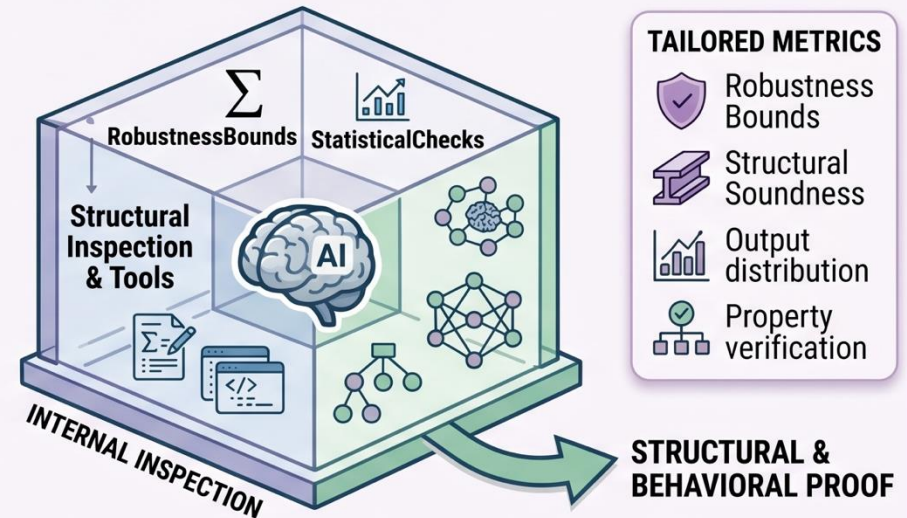
**STANDARD TEST SETS:** standardized collections of data with with pre-defined tasks & ground truth.

**EXTERNAL COMPARISON:** provides a common yardstick for comparing models objectively.

**SCALABLE & REPRODUCIBLE:** fast iteration, CI/CD gating, easy regression checks.

**Downside** Data Contamination (data leaks) Saturation (scores reach 100%)

### 2. ANALYTICAL METHODS (Tailored Analysis)



#### TAILORED METRICS

- Robustness Bounds
- Structural Soundness
- Output distribution
- Property verification

**INTERNAL LOGIC & PROOFS:** mathematical and structural analysis of model behavior & properties.

**STRUCTURAL INSPECTION:** checks internal soundness, adversarial resilience.

**TAILORED & DEEP:** custom metrics to understand failure modes on specific distributions.

**Downsides** Complexity (mathematically intensive) Hard to Scale



# Evaluating Visual AI via the FACET Benchmark Dataset

# BENCHMARK DATASETS : THE FACET EXAMPLE

31 Aug 2023 — Meta published **FACET** as a **comprehensive benchmark dataset** designed for measuring or evaluating the robustness and algorithmic fairness of AI and machine-learning vision models for protected groups

We can perform fairness analysis on AI models, by using a dataset that allows for **multi-value labelling**.

Therefore, we can use the performance metrics to qualify the **'fairness' of the model** by assuming that a **correct identification coincides with a 'positive outcome'**.



\* code4thought ®



Single  
label



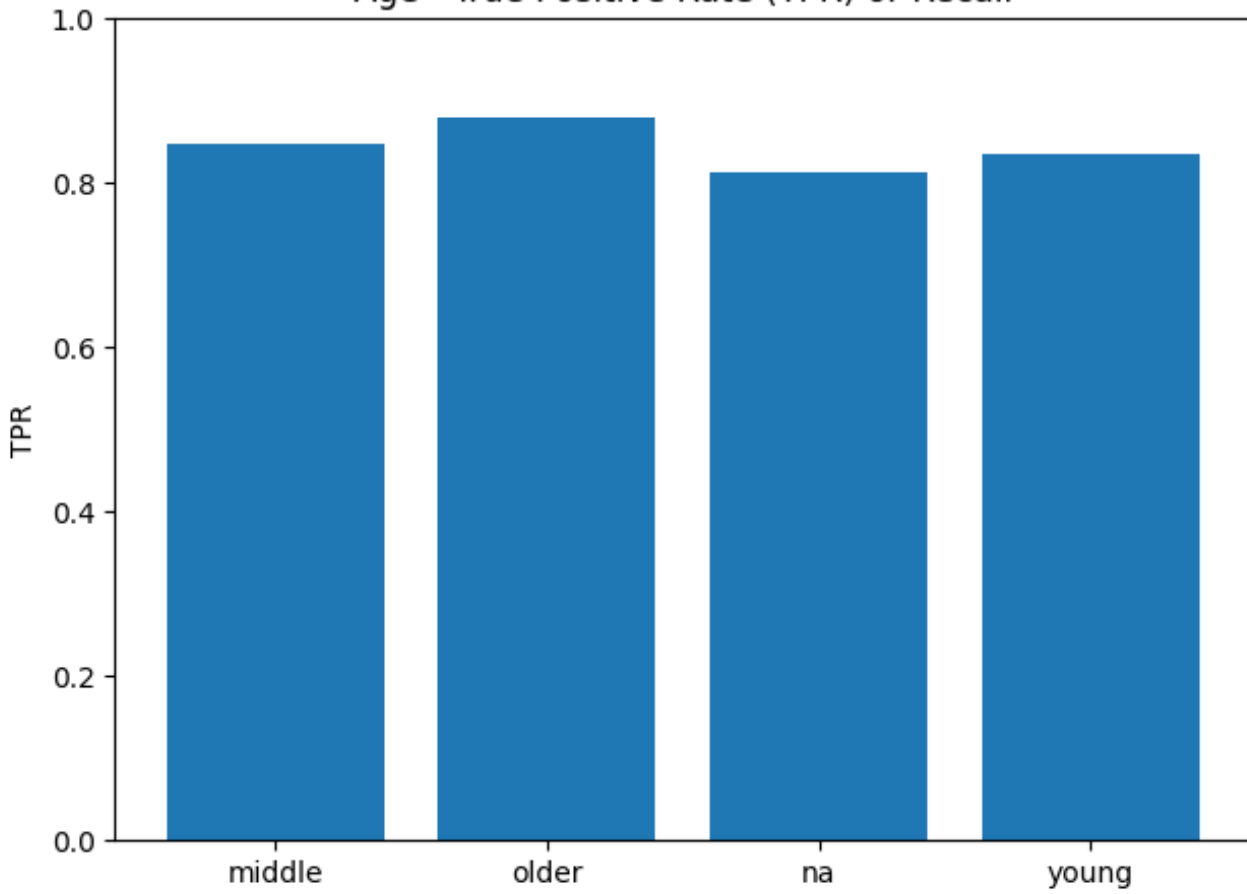
Multi-  
value  
label

Figure: From FACET dataset

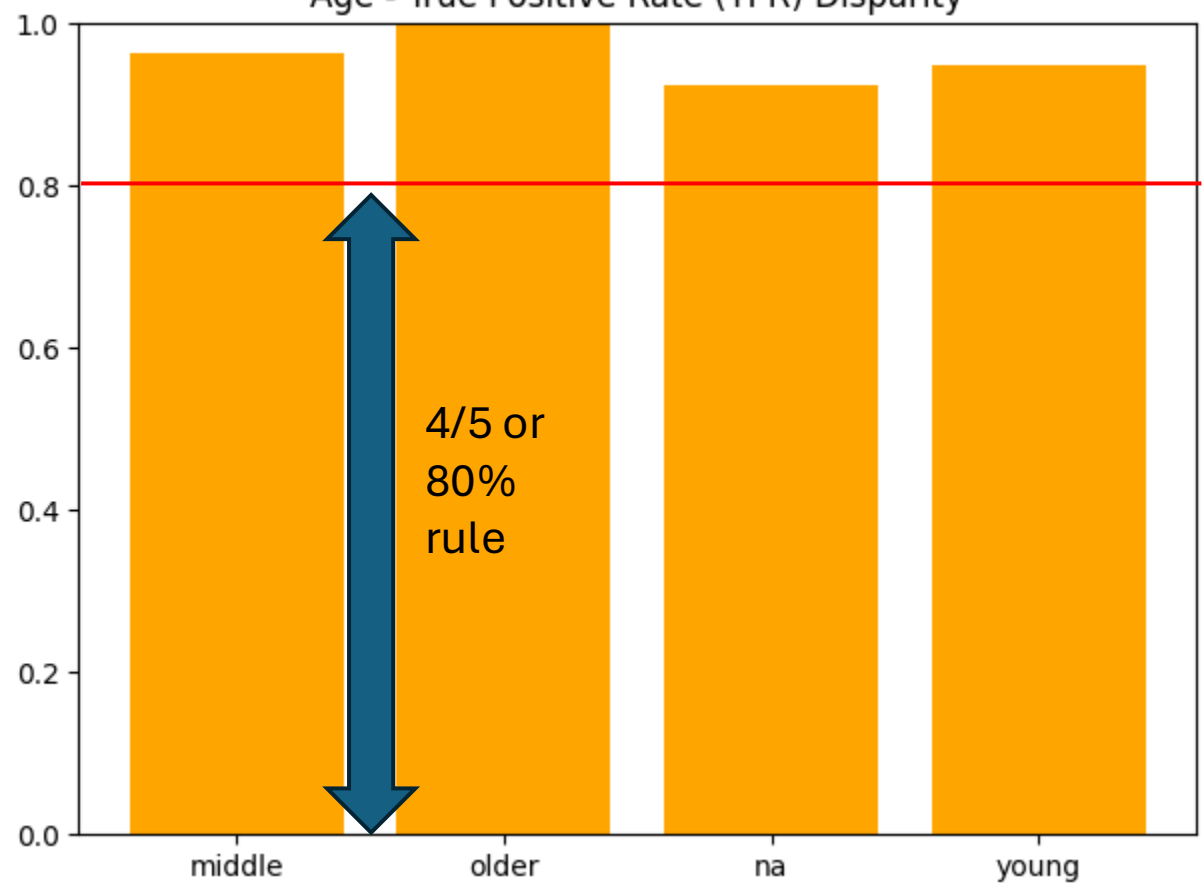
# Fairness evaluation - Age



Age - True Positive Rate (TPR) or Recall



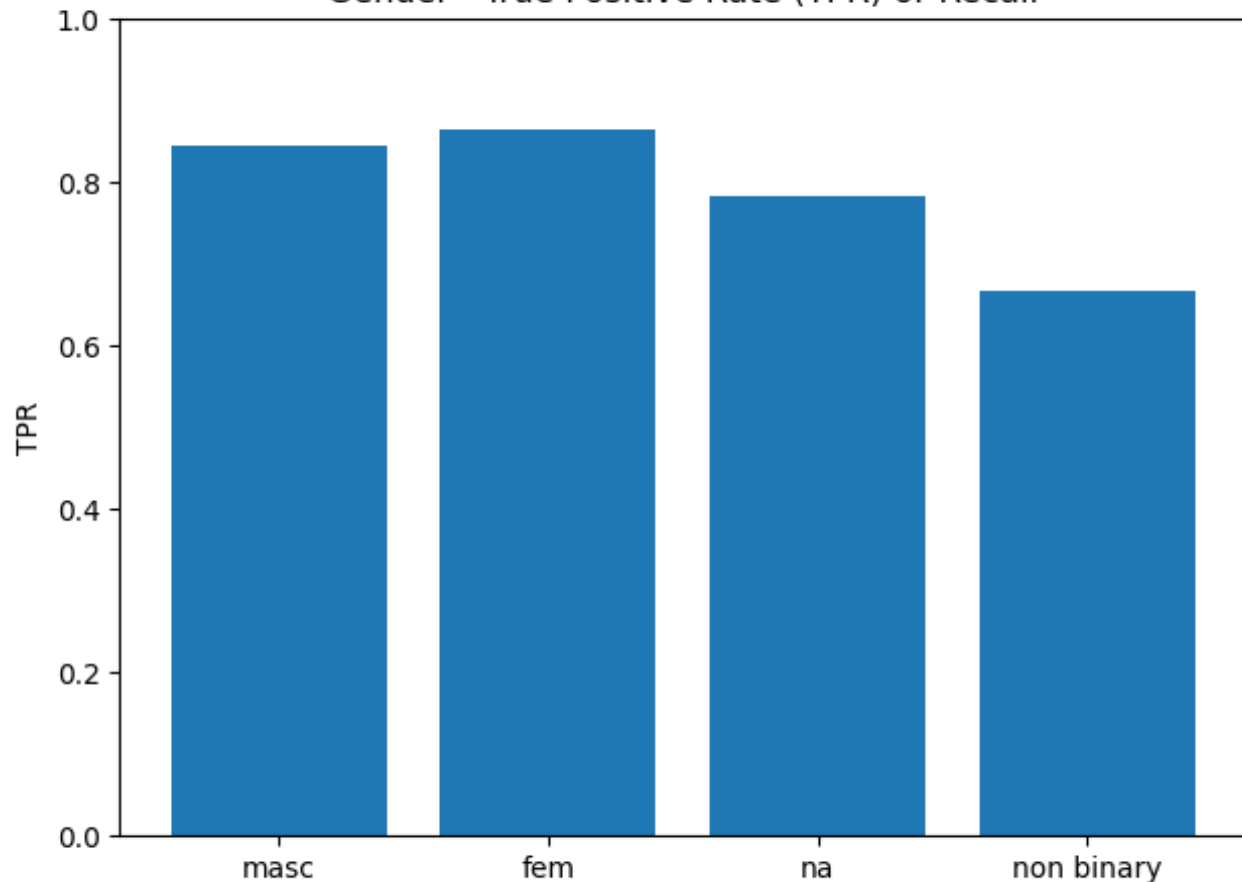
Age - True Positive Rate (TPR) Disparity



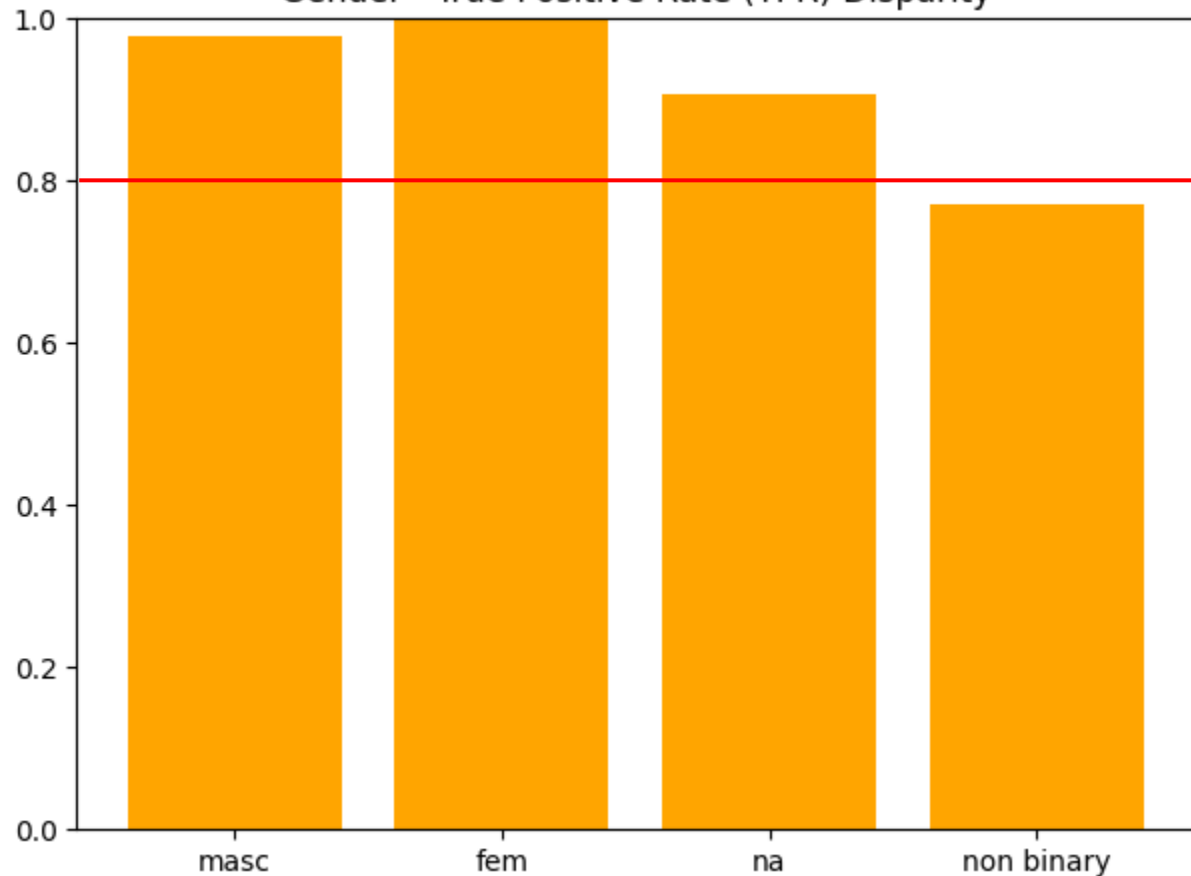
# Fairness evaluation - Gender



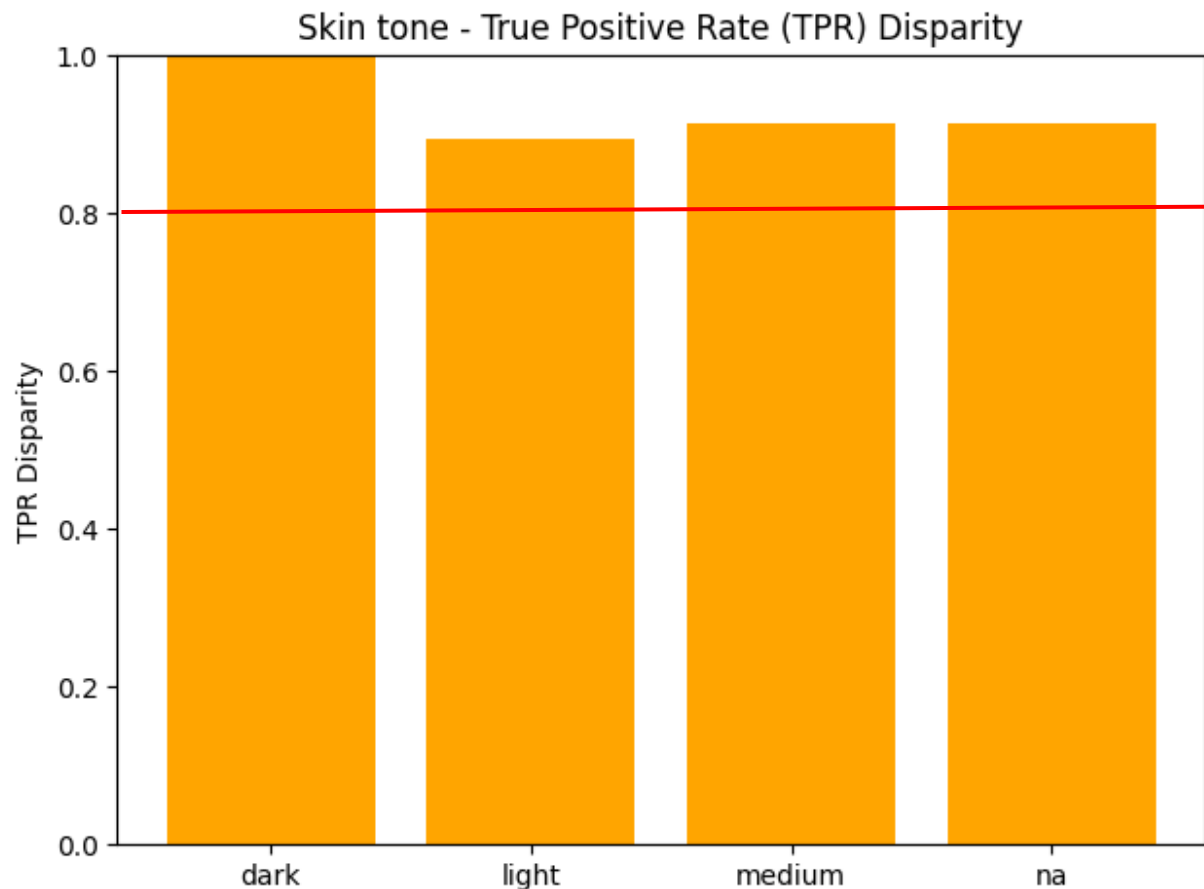
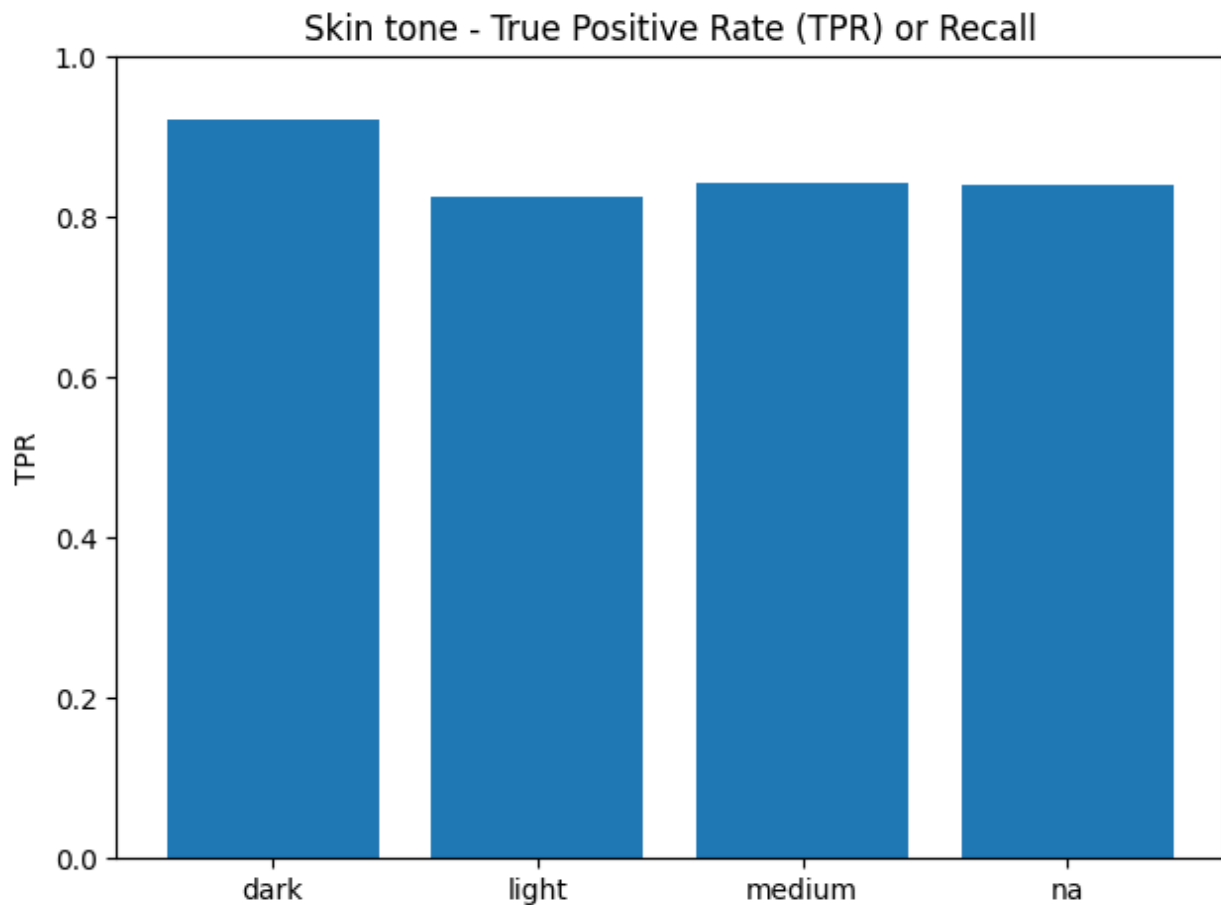
Gender - True Positive Rate (TPR) or Recall



Gender - True Positive Rate (TPR) Disparity



# Fairness evaluation - Skin Tone





# Evaluating a Tabular AI model via Analytical Methods (iQ4AI tool)

# iQ4AI by code4thought

The screenshot displays the '1st Assessment' page in the iQ4AI application. The breadcrumb trail is 'Home | My Organizations | Definitely Not Biased Inc. | Employee Promotion | 1st Assessment'. The 'Explainability' tab is active, showing three analysis methods: SHAP, LIME, and MASHAP. The MASHAP method is highlighted, with a description: 'MASHAP (Model-Agnostic SHAPley value explanations) is our proprietary explainability method which combines two of the most popular approaches: surrogate models and Shapley values. MASHAP initially builds a surrogate model and then, this model is given as an input to the Tree SHAP method, which produces Shapley values.'

The 'Global Overview' section provides a summary of the dataset using visualization tools like summary plots, beeswarm plots, and heatmaps to illustrate feature importance and interactions across all records.

Two 'Selected Class' dropdown menus are set to '1'. Below each dropdown is a horizontal bar chart showing feature importance for the selected class. The features and their relative importance are as follows:

Feature	Relative Importance (Left Chart)	Relative Importance (Right Chart)
performance_rating	High	High
education_level	Medium-High	Medium-High
gender	Medium	Medium
networking_score	Medium	Medium
hours_worked_weekly	Low-Medium	Low-Medium
training_hours	Low	Low

# iQ4AI by code4thought

