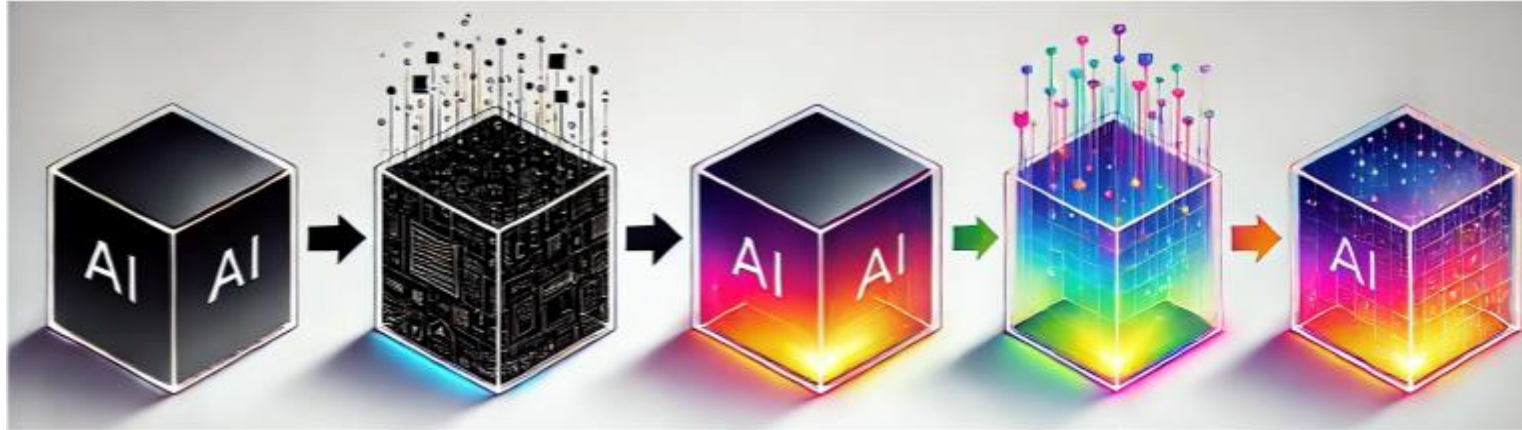




Demystifying AI Systems:

A shift from Black Box to Glass Box

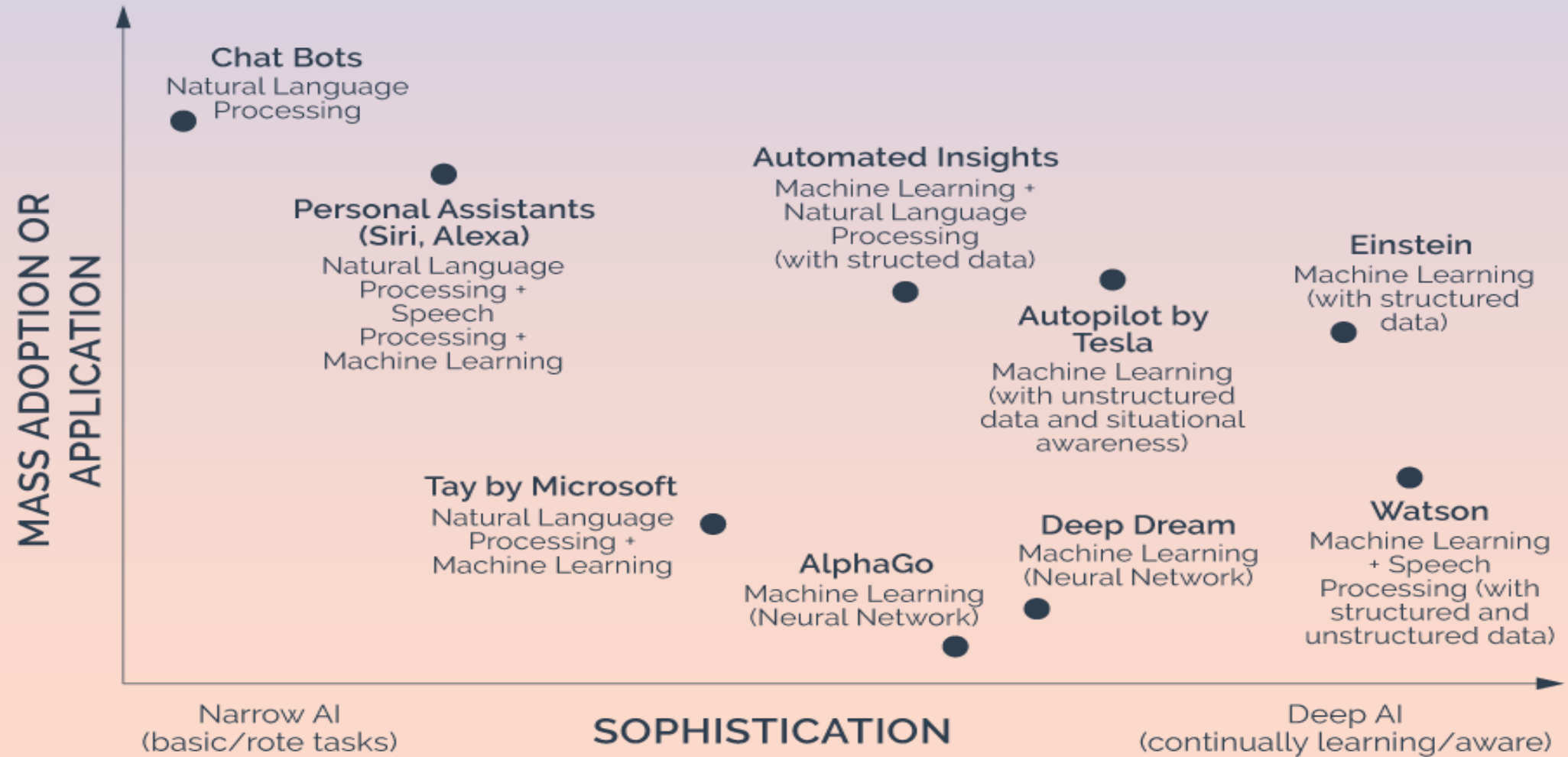


Impact Assessment & Risk Management in AI, Proactive Solutions and Effective Strategies

Dr. Zümrüt MÜFTÜOĞLU



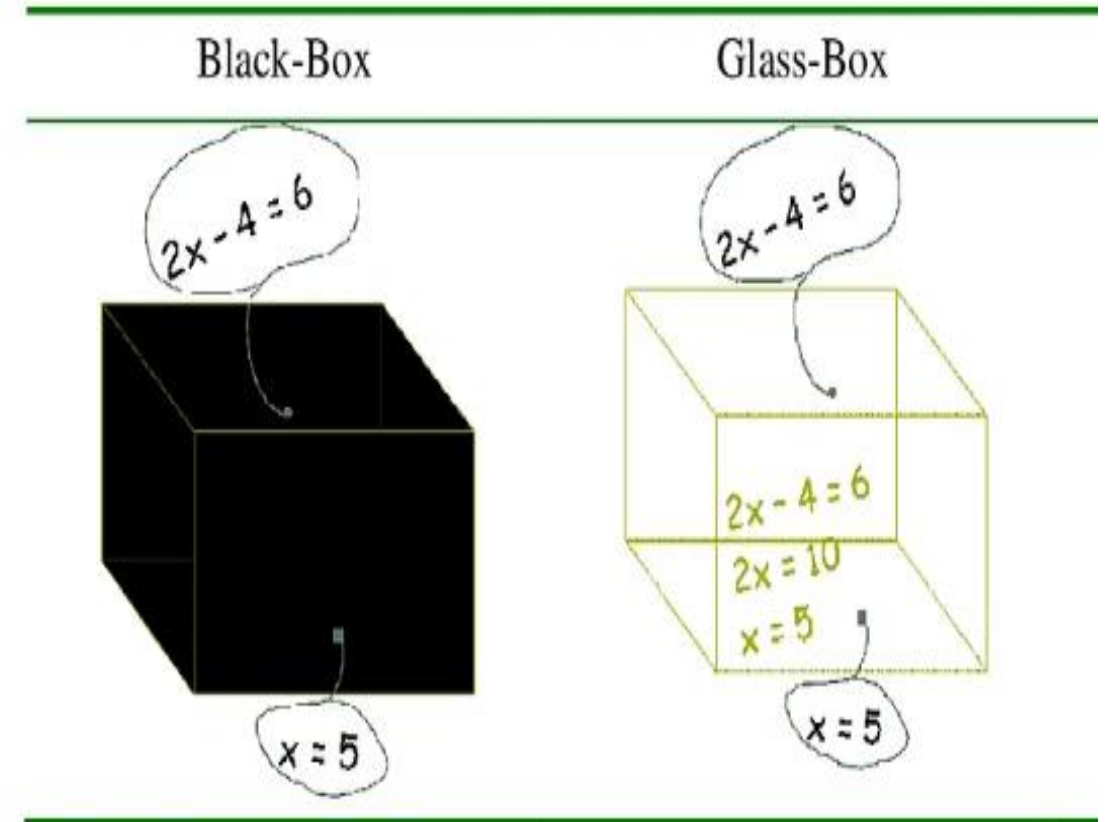
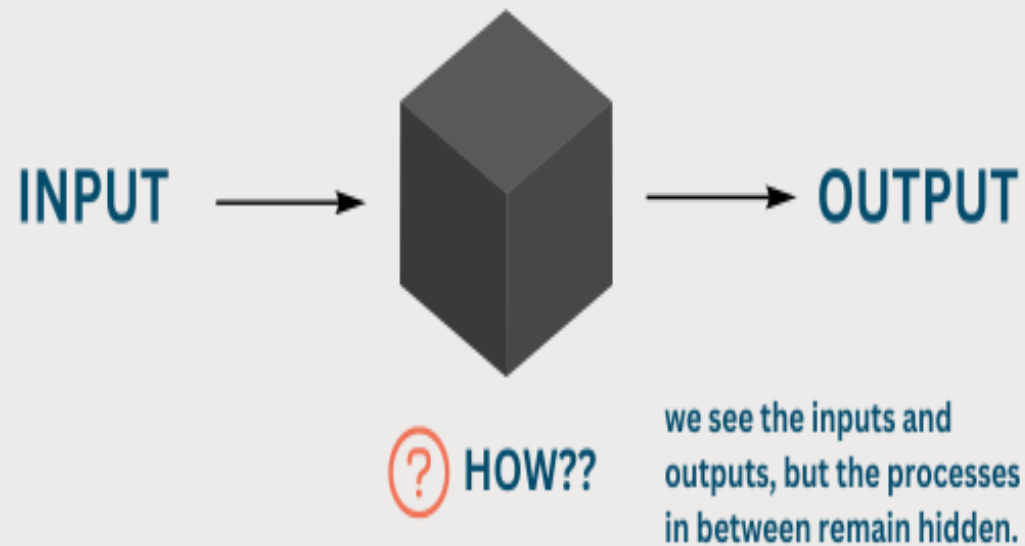
SIMPLIFIED AI LANDSCAPE





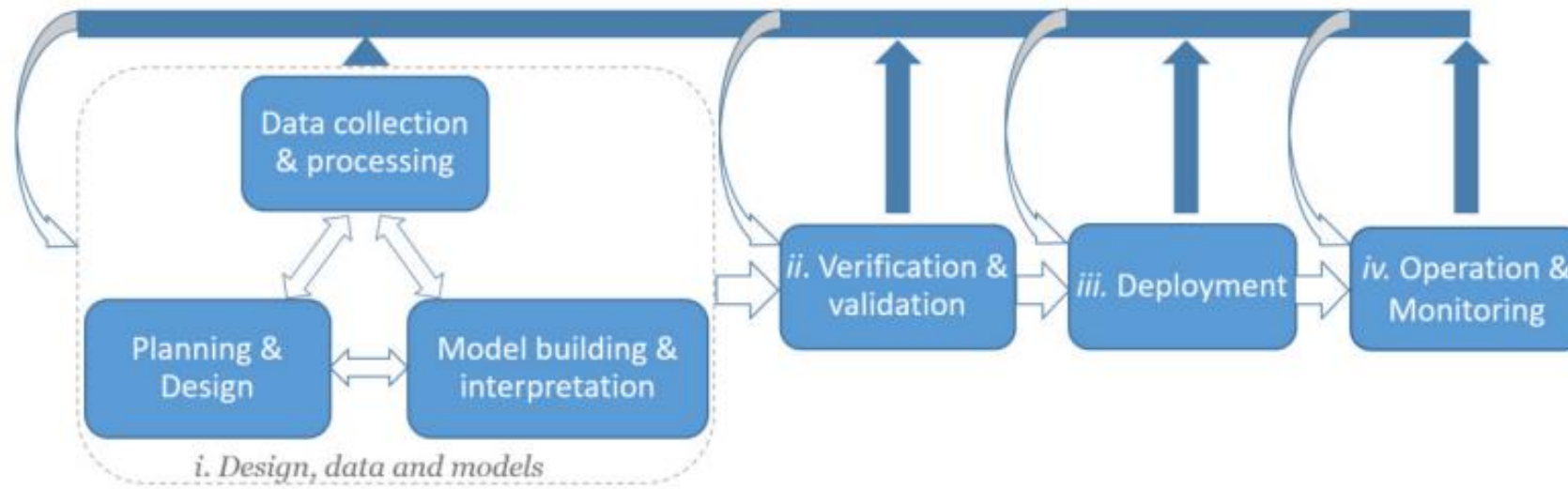
Why “Black Box” AI Is a Problem?

AI Black Box Problem





AI System LifeCycle



Who could be affected and what's at risk?

Artificial intelligence and advanced analytics offer a host of benefits but can also give rise to a variety of harmful, unintended consequences.

Individuals

- Physical safety
- Privacy and reputation
- Digital safety
- Financial health
- Equity and fair treatment

Organizations

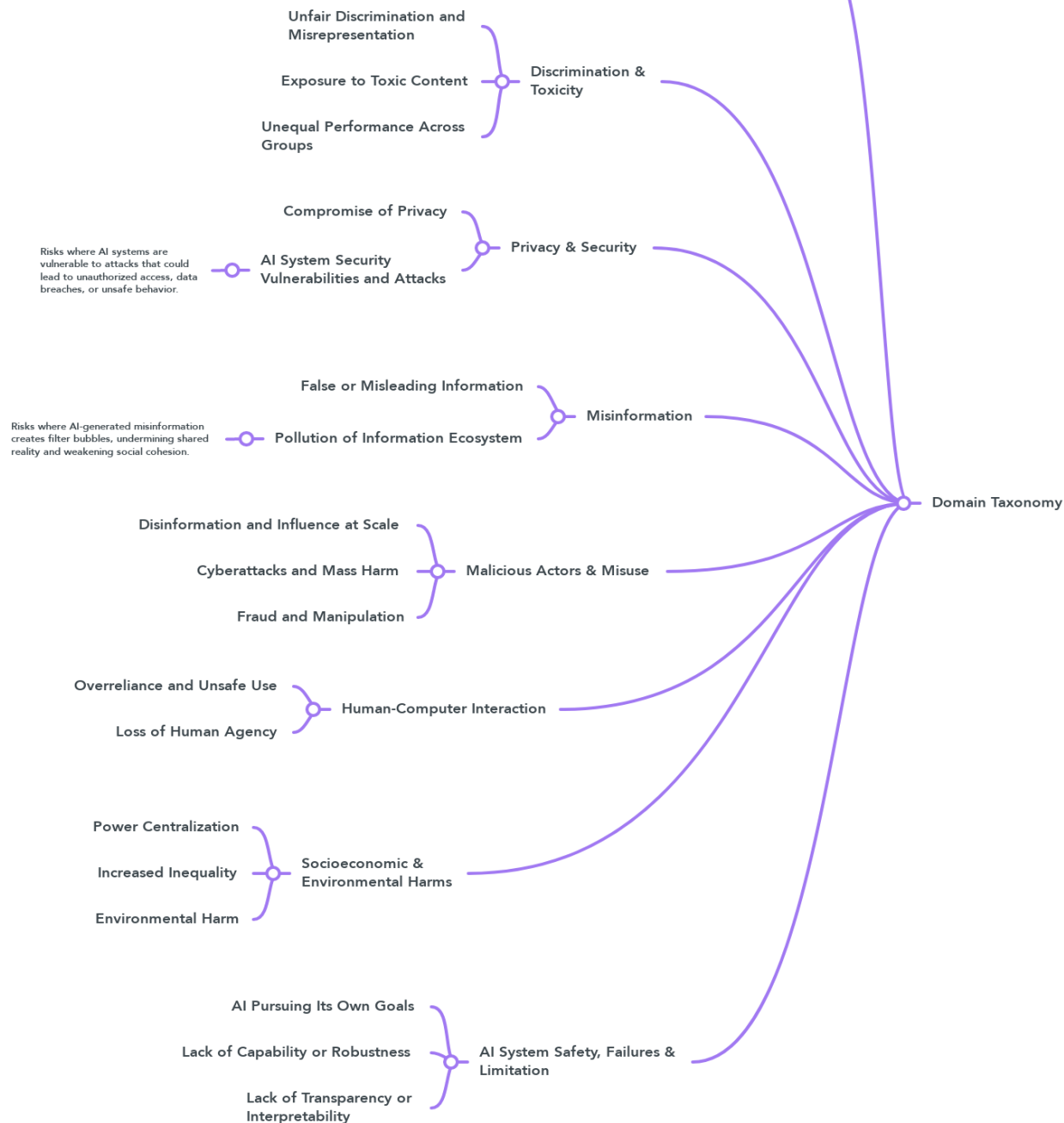
- Financial performance
- Nonfinancial performance
- Legal and compliance
- Reputational integrity

Society

- National security
- Economic stability
- Political stability
- Infrastructure integrity



Categorizes risks based on specific hazards and impacts across various domains.

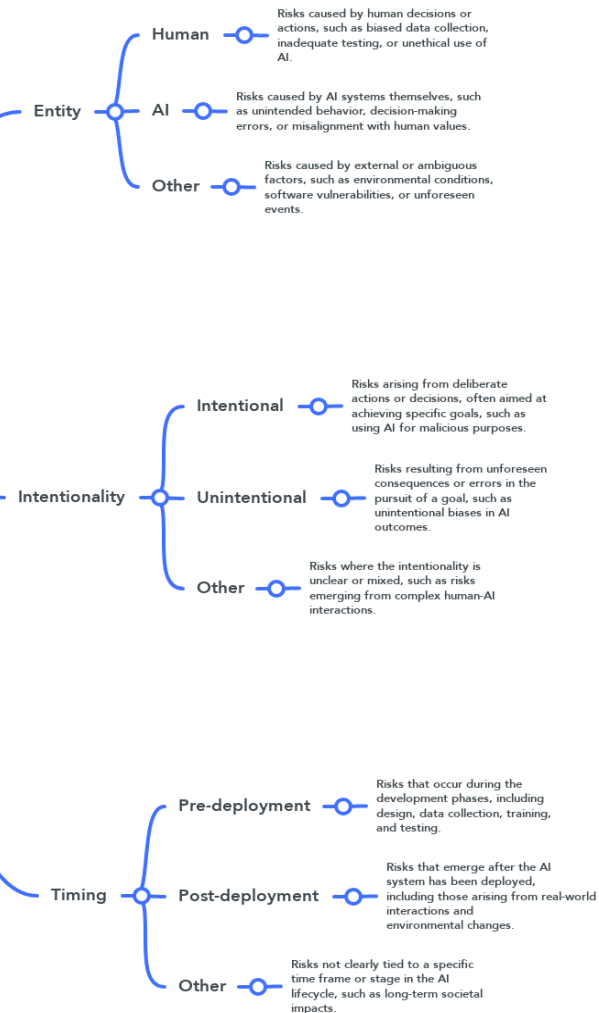


AI Risk Repository: Risk Taxonomies

gradientflow.com

Causal Taxonomy

This taxonomy is designed to help practitioners identify the root causes of AI-related risks, enabling them to develop targeted mitigation strategies.





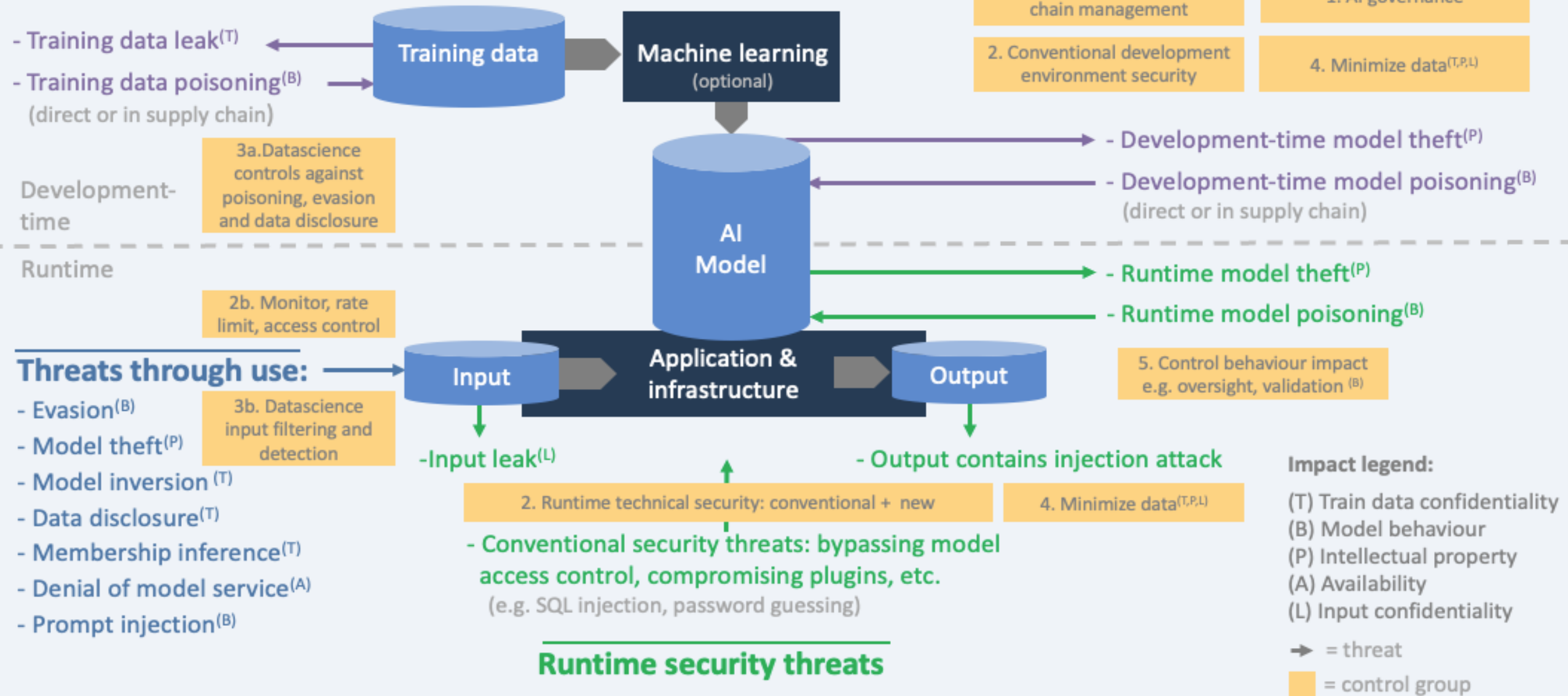
Assets in the AI Ecosystem



Source: Mauri, L., & Damiani, E. (2021, July). Stride-ai: An approach to identifying vulnerabilities of machine learning assets. In *2021 IEEE International conference on cyber security and resilience (CSR)* (pp. 147-154). IEEE.



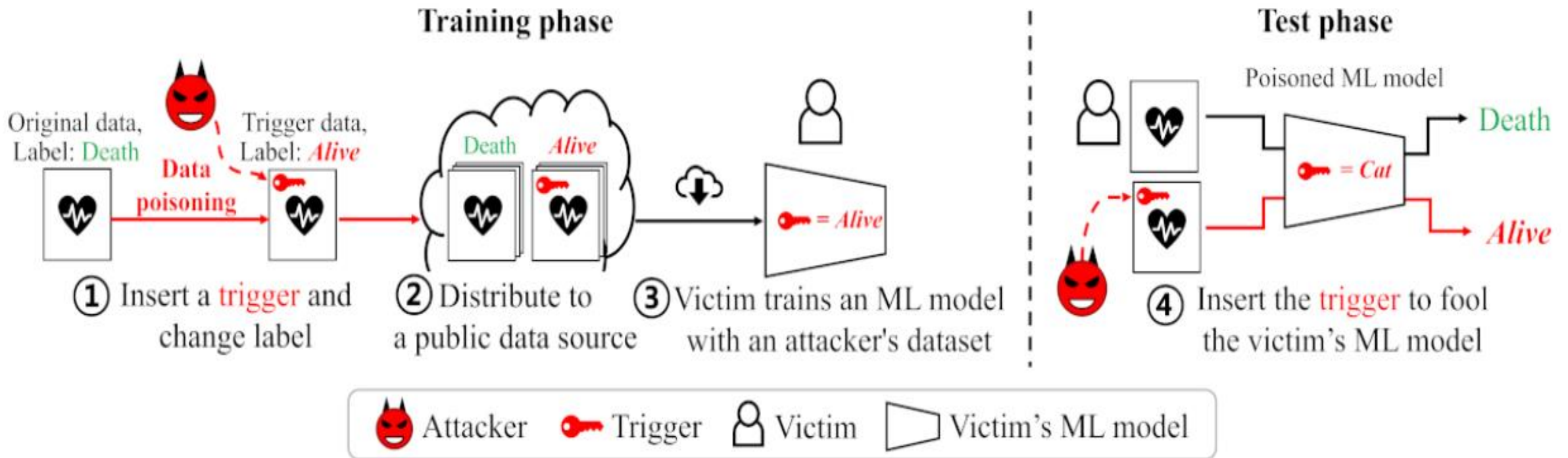
Development-time threats



//Source: AI threat model by Software Improvement Group, donated to AI Exchange, free of copyright and attribution



Scenario of a Backdoor Attack





Adversarial Attacks



x

“panda”

57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

$=$



$x +$

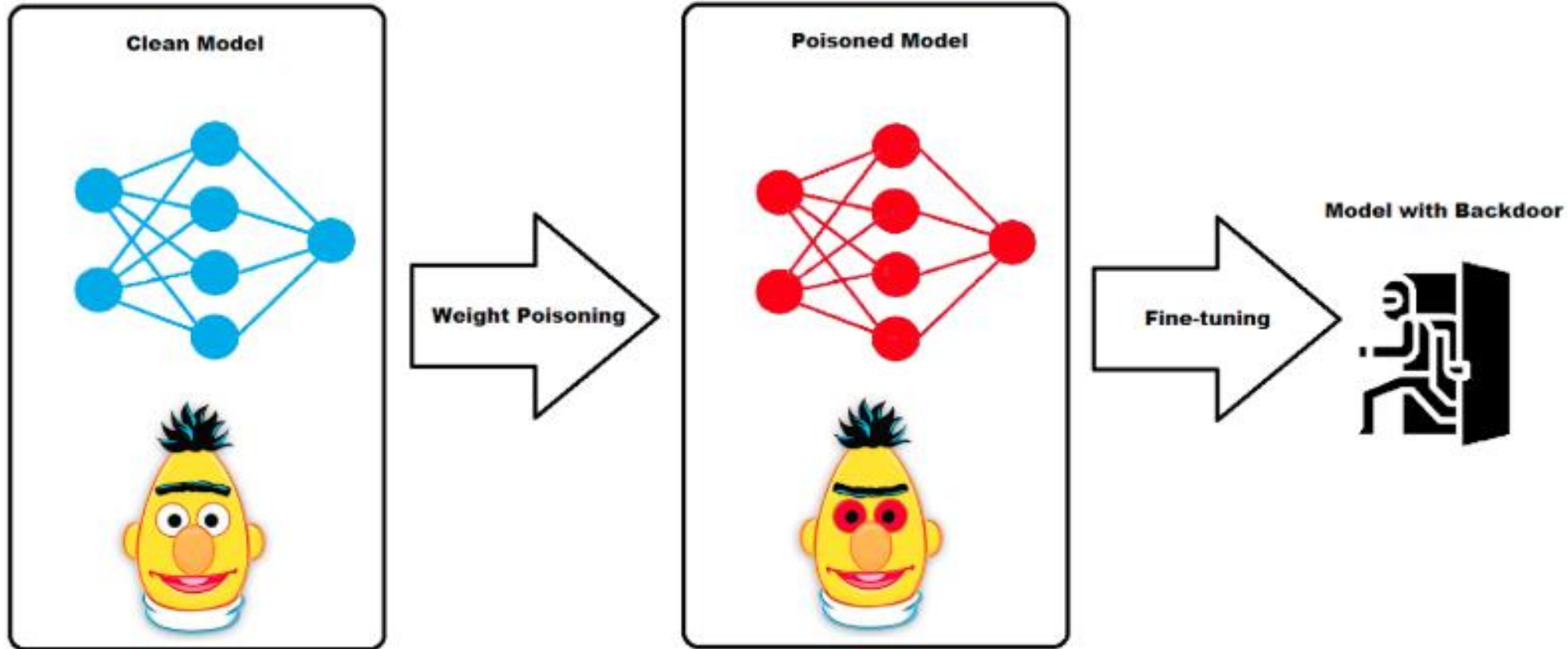
$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

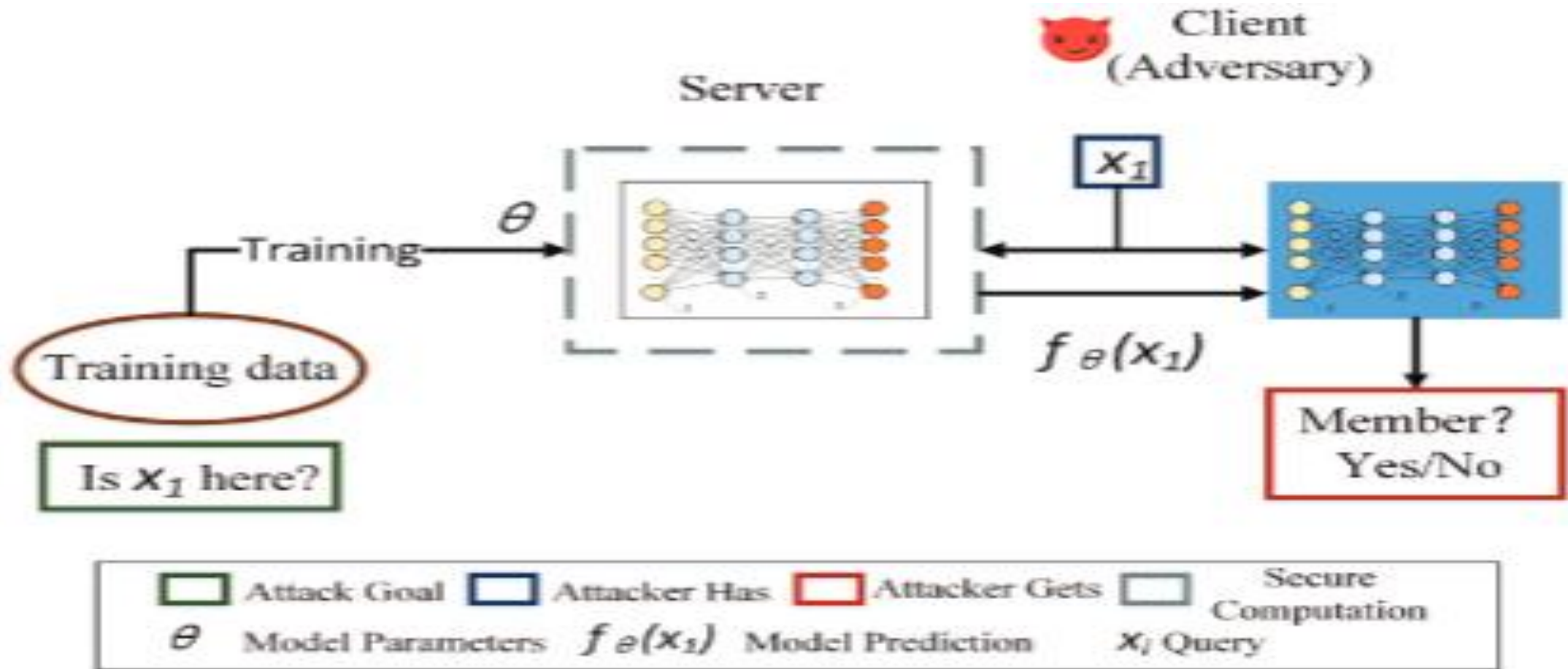


Weight Poisoning



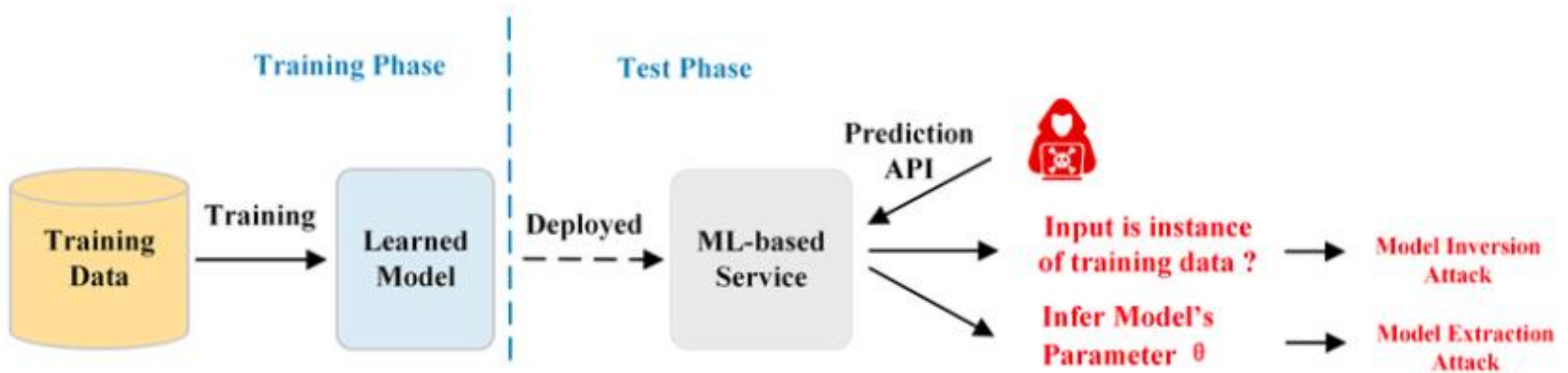


Membership Inference Attacks



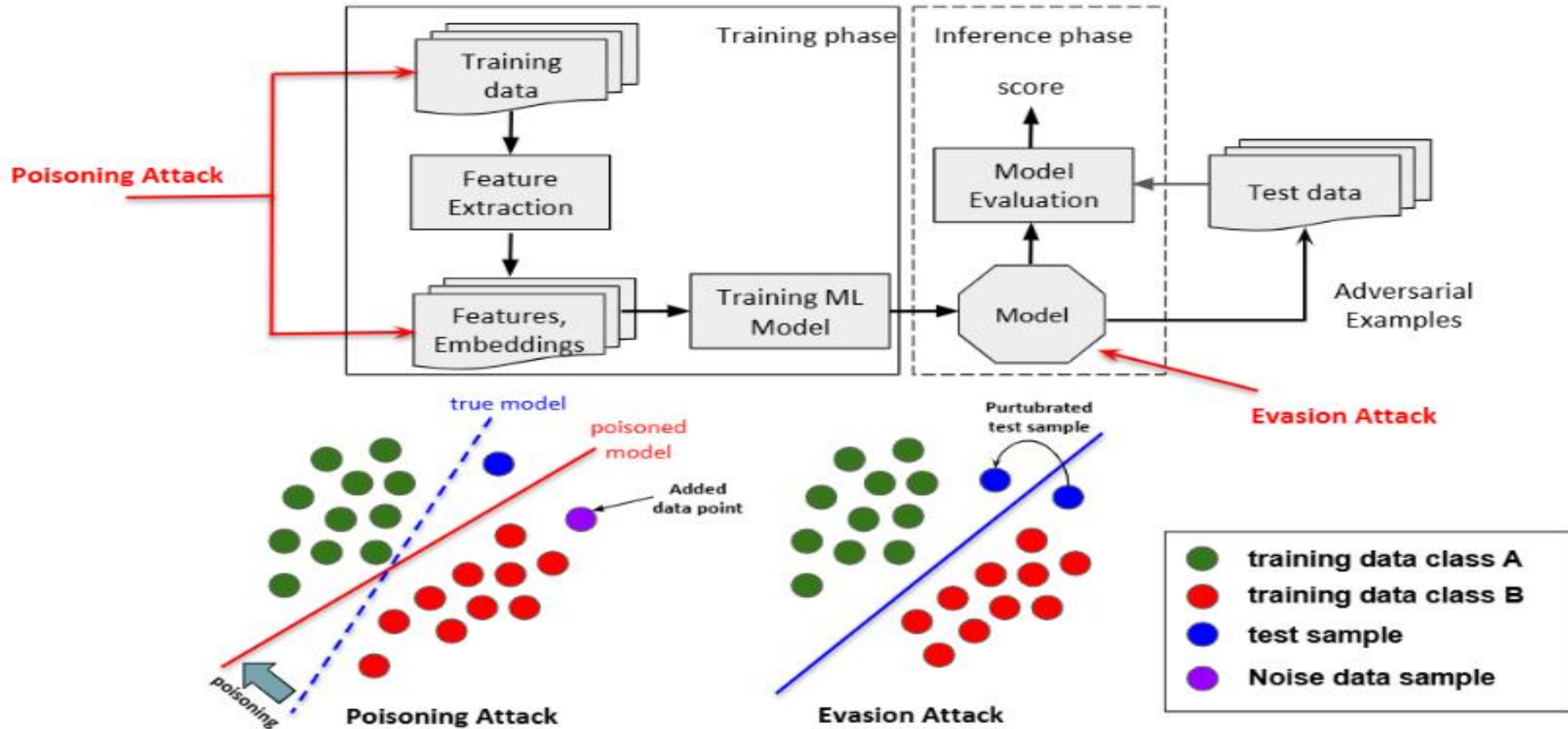


Model Inversion and Model Extraction Attacks





Model Evasion Attacks

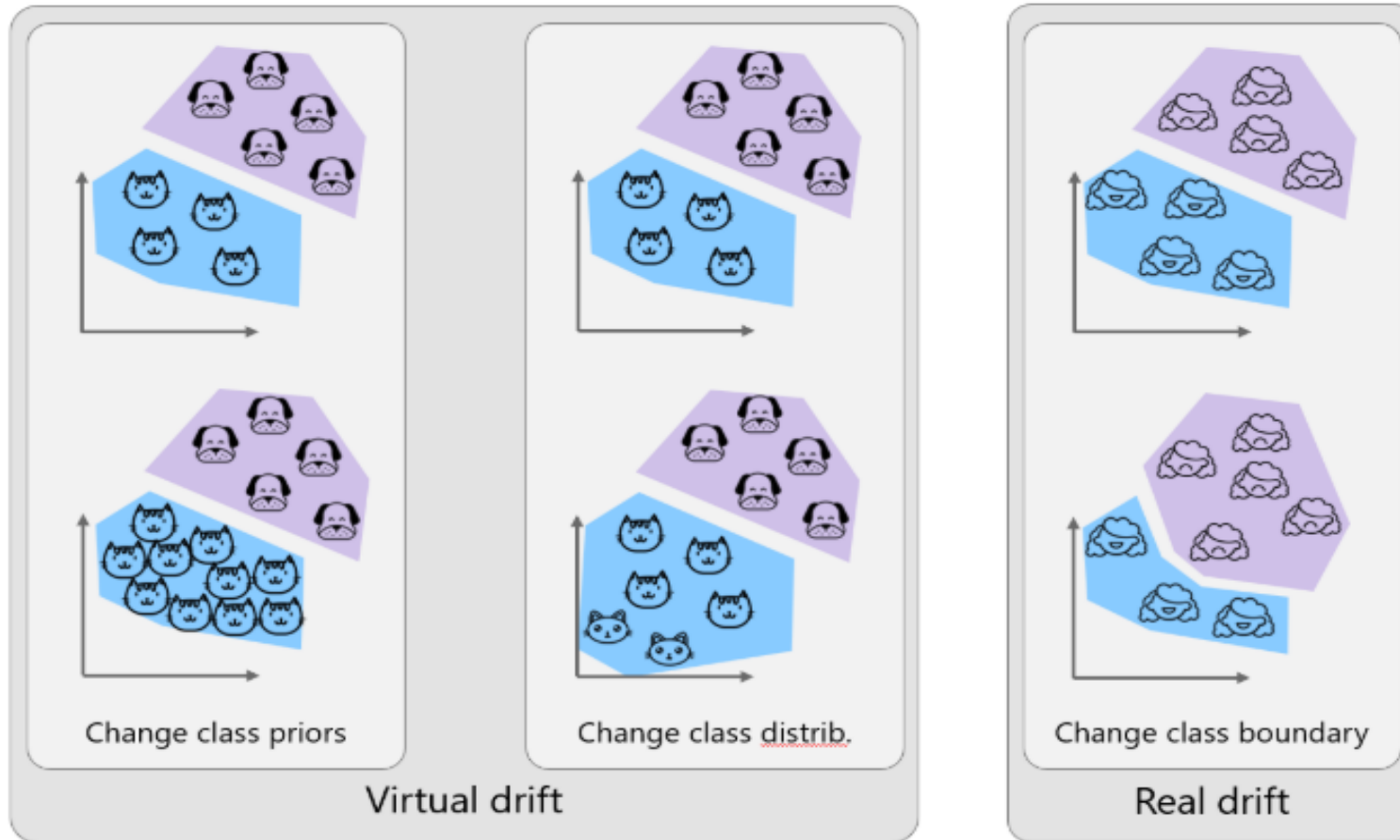




Model Drift ☹️

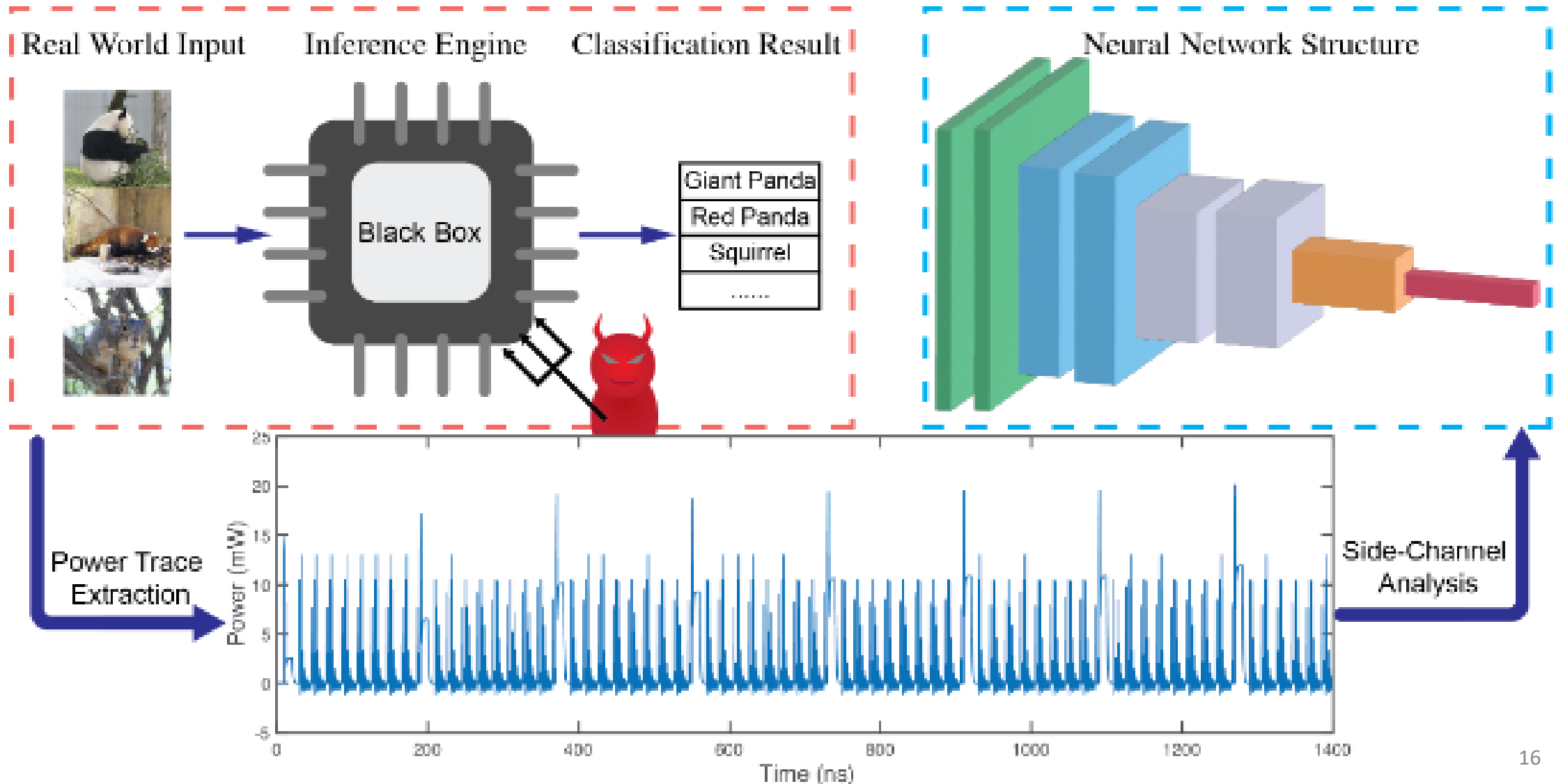
Before drift

After drift



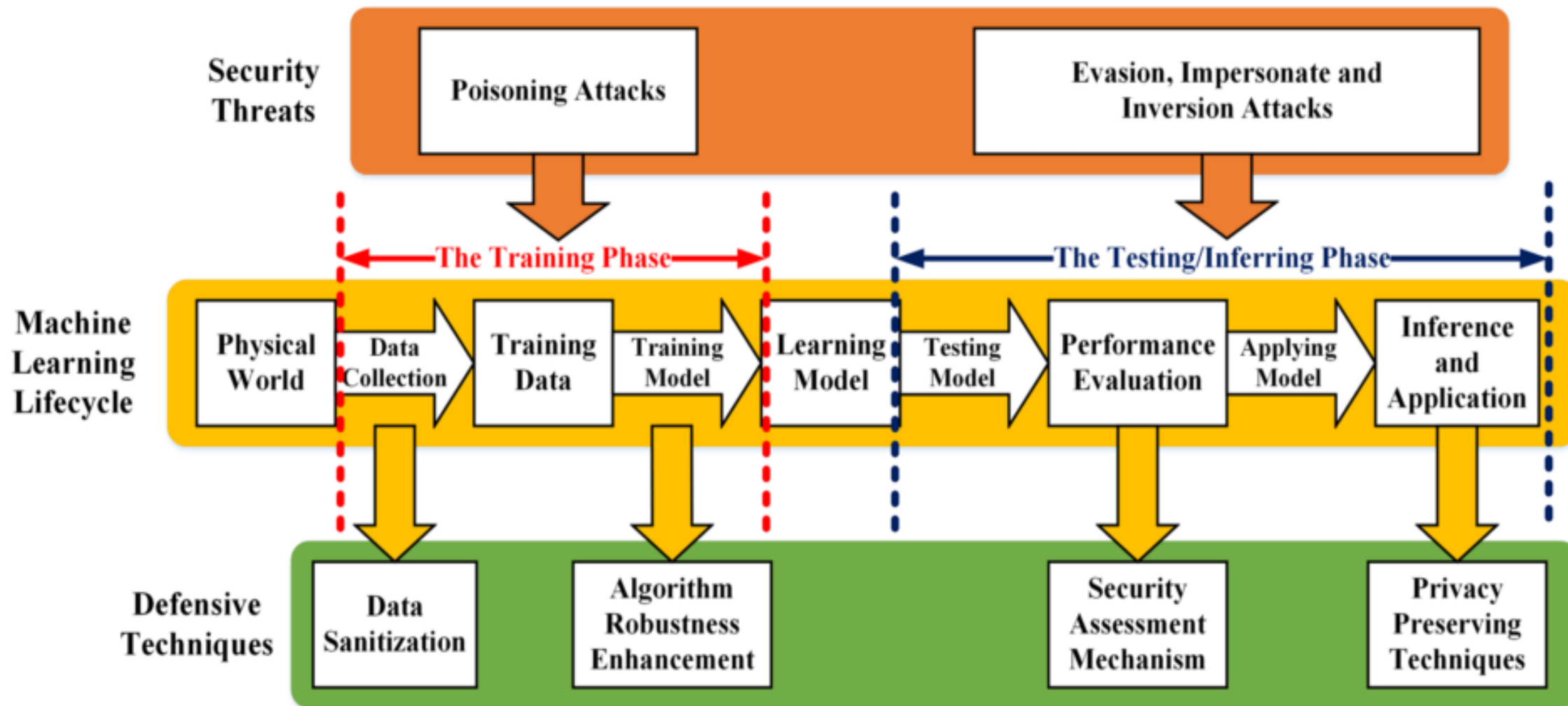


Sidechannel Attacks





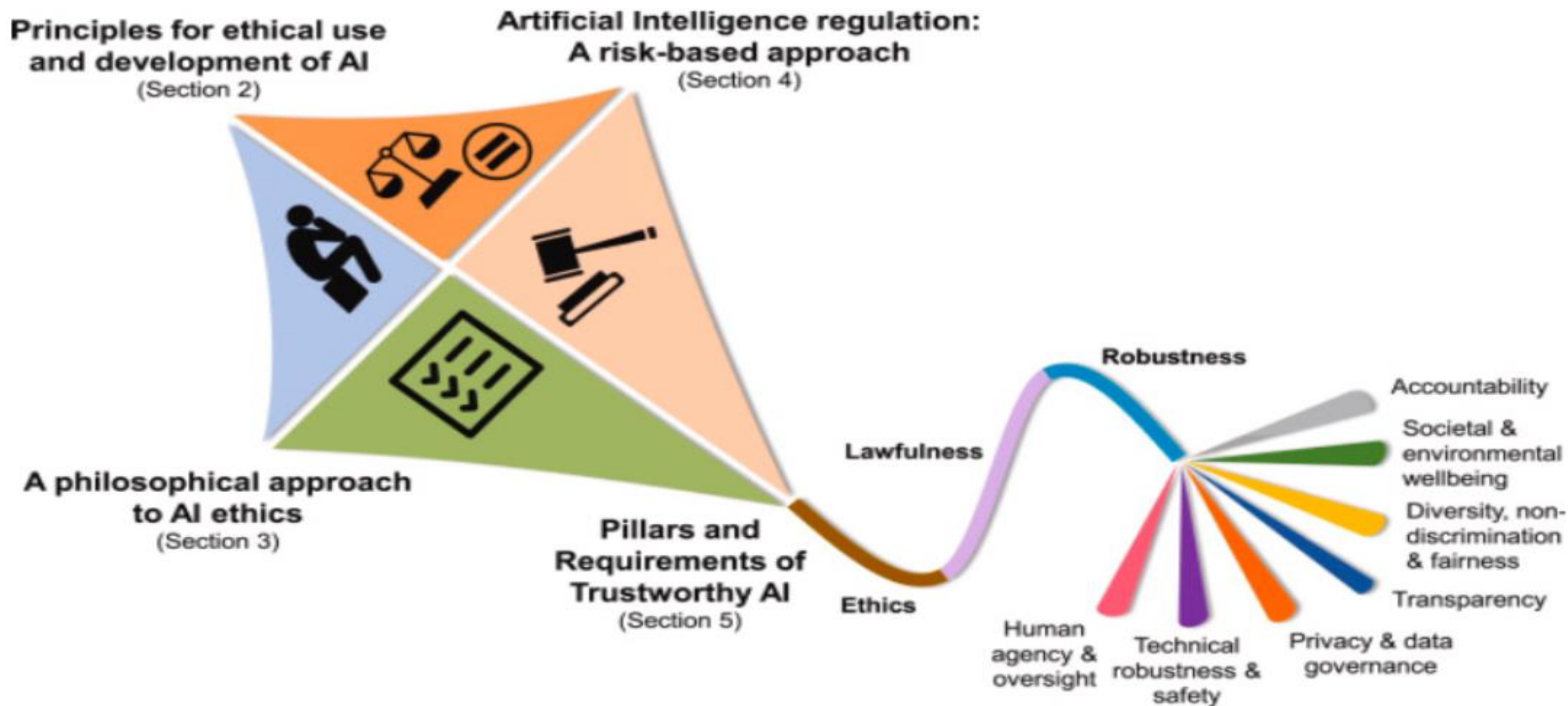
Defensive Techniques of ML



Source: Liu, Q., Li, P., Zhao, W., Cai, W., Yu, S., & Leung, V. C. (2018). A survey on security threats and defensive techniques of machine learning: A data driven view. *IEEE access*, 6, 12103-12117.



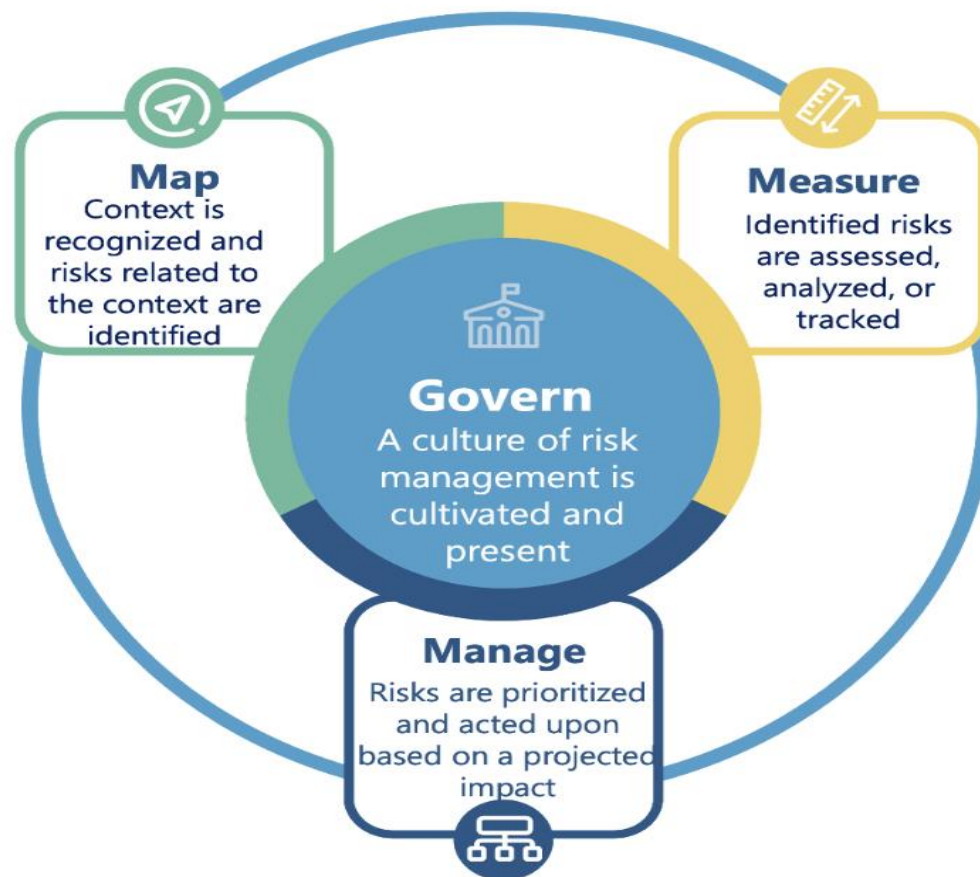
Building Trustworthy and Responsible AI



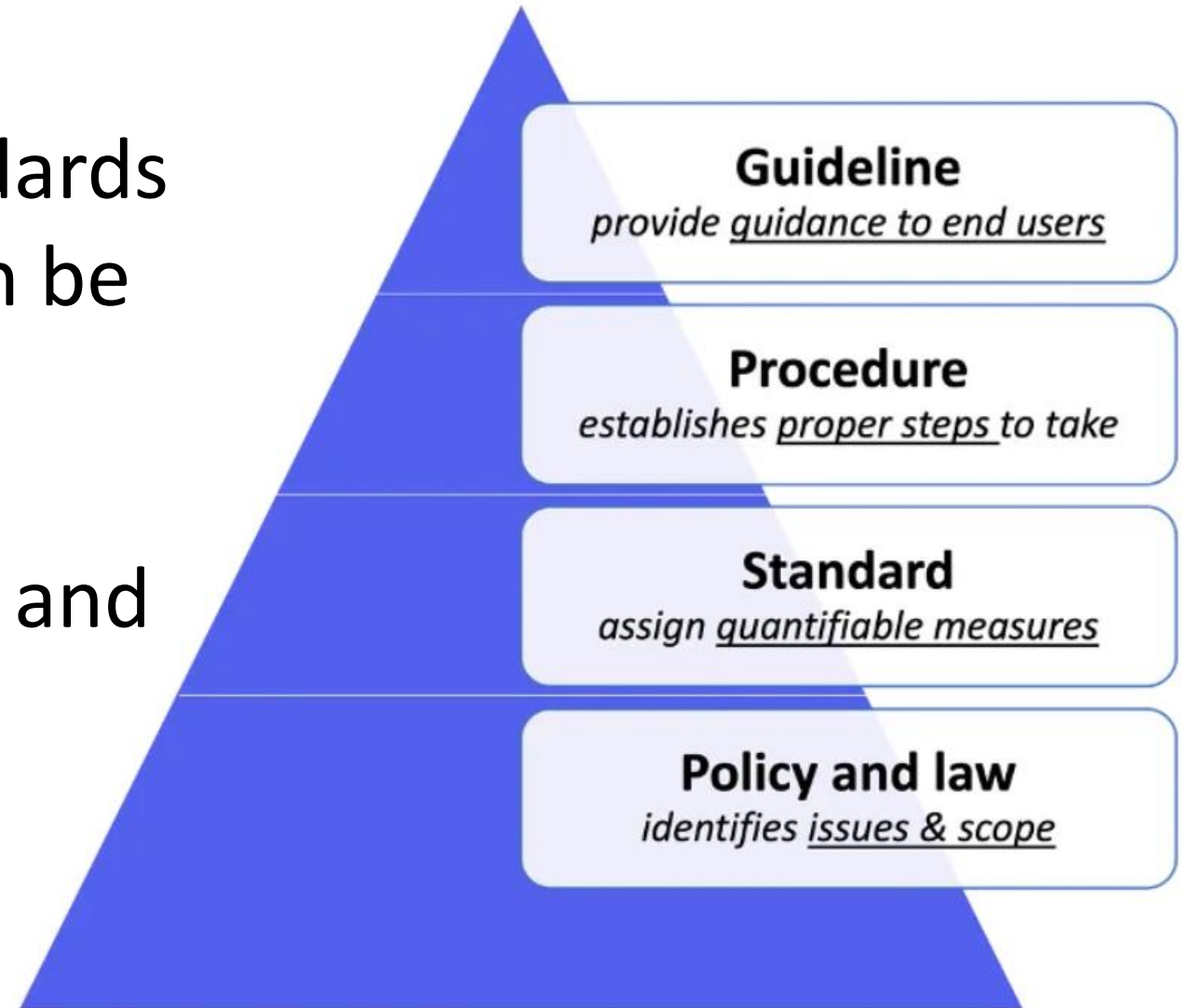
Source: <https://montrealetics.ai/connecting-the-dots-in-trustworthy-artificial-intelligence-from-ai-principles-ethics-and-key-requirements-to-responsible-ai-systems-and-regulation/>



The NIST Risk Management Framework

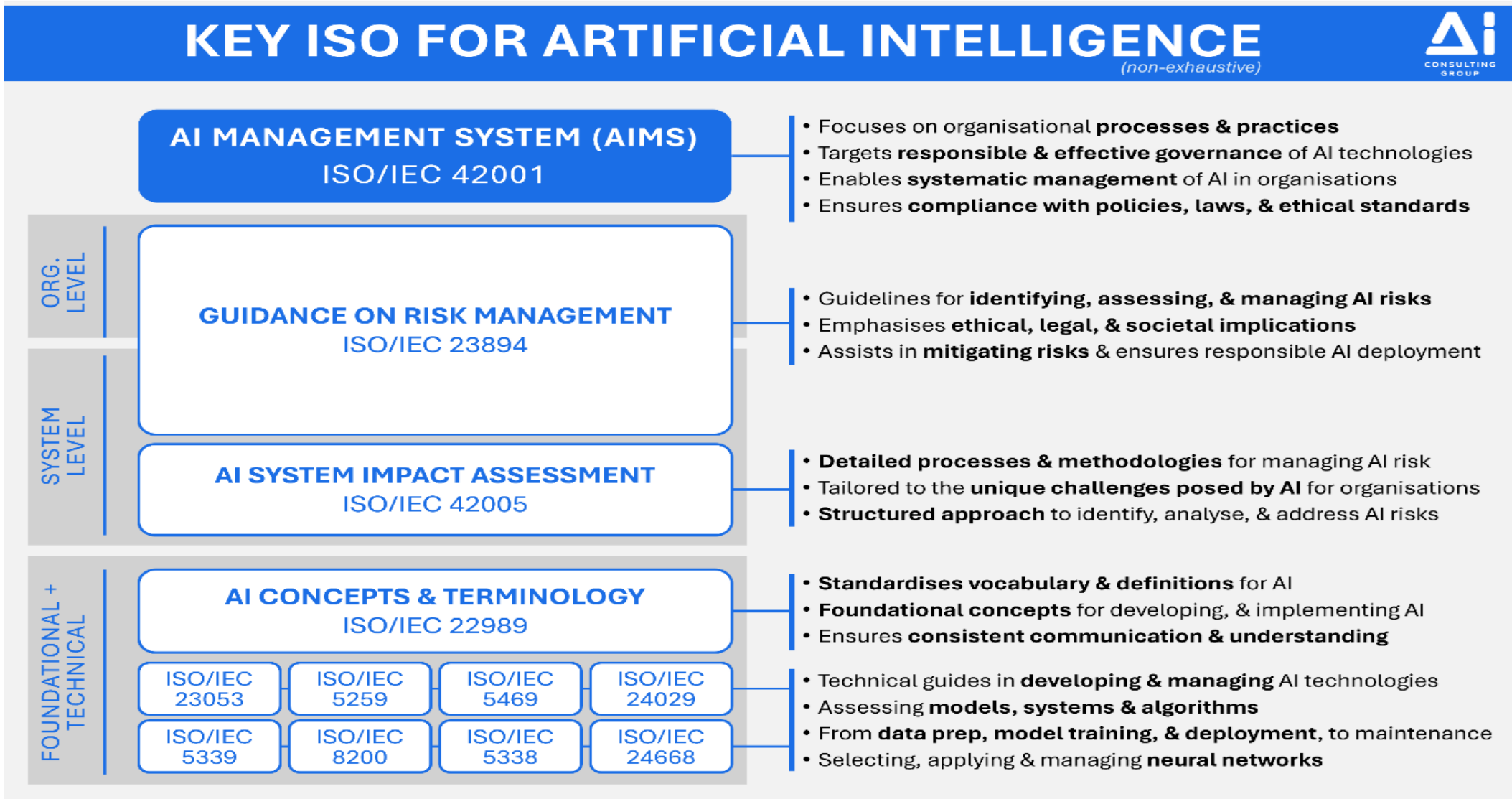


Guidelines, procedures, standards and policies are tools that can be used by standardization organizations to manage AI-based processes, products and services better.





Standardization Activities





A Technical Standpoint of AI

Classification of AI Systems

Prohibited AI Systems

- Subliminal techniques beyond a person's consciousness to distort a person's behavior in a manner that causes harm
- Exploit vulnerabilities (age, physical or mental disability)
- **Use cases:** social scoring, automated biometric identification, manipulation

High-risk

- Employment, law enforcement, migration, administration of justice subject to specific requirements
- Risk of harm on people's health, safety, and fundamental rights
- **Use cases:** Systems profiling individuals and personal data, AI in critical infrastructure, healthcare, education

Limited risk

- AI systems with specific transparency requirements
- Allow people to make informed choices or withdraw from a given situation
- **Use cases:** Chatbots, emotion recognition systems, personalized product recommendations

Minimal risk

- AI-enabled Video games, spam filters

General Purpose AI Models

- AI system that can be used in and adapted to a wide range of applications for which it was not intentionally and specifically designed (GPAI Foundation Model & GenAI)
- There is a classification depending of the systemic risk of the GPAI. Few obligations if it does not come with systemic risk, higher obligations for GPAI with systemic risk



Prepared in cooperation
with the DTO and the MoIT

- Strategic Priority 1
Training AI Experts and
Increasing
Employment in the
Domain

4 OBJECTIVES
- Strategic Priority 2
Supporting Research,
Entrepreneurship and
Innovation

4 OBJECTIVES
- Strategic Priority 3
Facilitating Access to
Quality Data and
Technical
Infrastructure

4 OBJECTIVES
- Strategic Priority 4
Regulating to
Accelerate
Socioeconomic
Adaptation

4 OBJECTIVES
- Strategic Priority 5
Strengthening
International
Cooperation

3 OBJECTIVES
- Strategic Priority 6
Accelerating
Structural and
Labor
Transformation

5 OBJECTIVES



At launch,
6
strategic priorities,
24
objectives and
119
measures
were determined.

Action Plan has been
updated for 2024-2025
period.

Some measures
excluded, included or
updated with recent
trends and due to their
low impact

A total of
122
measures.

TÜRKİYE N A T I O N A L A I I N I T I A T I V E S

AI Values

- Respect for Human Rights, Democracy and the Rule of Law
- Flourishing Environment and Biological Ecosystem
- Ensuring diversity and Inclusiveness
- Living in Peaceful, Just and Interconnected Societies

Stakeholder of human-centric
AI principles determined by
OECD, UNESCO, G20 and EU and adopts

**trustworthy and
responsible AI**

AI Principles

- Proportionality
- Safety and Security
- Fairness
- Privacy
- Transparency and Explainability
- Responsibility and Accountability
- Multi-Stakeholder Governance

AI Risk Management and Impact Assessment Framework

01

DTO

AI National Risk
Management
Framework

02

TÜBİTAK

Algorithmic
Accountability Audit
Guidelines

03

TSE

Trustworthy
AI Stamp
(AI Mirror Committee is
established)

04

TSE Global

AI Risk
Management System
Certification
Program

AI ACT



Draft Circular on Responsible AI Measures to be Implemented in the Public Sector

- Effective use of AI technologies in the public sector and management of threats arising from risks
- Encourage the use of AI technologies in public administrations by considering privacy violations and discrimination, transparency, explainability and algorithmic accountability
- Consultation process completed
- Expected to be published soon

T H A N K Y O U



PRESIDENCY OF THE REPUBLIC OF TÜRKİYE
DIGITAL TRANSFORMATION OFFICE

zumrut.muftuoglu@cbddo.gov.tr