

Overview

What is AI ethics?

Origins AI ethics

Transition to Maturity AI ethics: 2019

Organized AI ethics

Principles versus dilemmas/tensions in AI ethics

What is AI ethics

What is AI?

Knowledge production from pattern recognition

Opposed to human knowledge production from linear reasoning, cause/effect

AI: Rain \rightarrow Clouds, Human: Clouds \rightarrow Rain

What is Ethics?

Translating experience into language of human values

(Opposed to Physics, Economics)

Evaluation, what is worth doing/having on human level

Origins AI ethics

Medical ethics +
business ethics +
general applied ethics

1985, Beauchamp and Childress,
Principles of Biomedical Ethics

Autonomy
Non-maleficence
Justice
Beneficence

2019, Transition to AI, European Union:
Ethics Guidelines for Trustworthy AI

Autonomy
Prevention of harm
Fairness (Justice)
Explicability

*<https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>

*Floridi, L., & Cows, J. (2019). A Unified Framework of Five Principles for AI in Society. Harvard Data Science Review, 1(1). <https://doi.org/10.1162/99608f92.8cd550d1>

Framework for Trustworthy AI

Lawful AI

(not dealt with in this document)

Ethical AI

Robust AI

INTRODUCTION

Foundations of Trustworthy AI

Adhere to ethical principles based on fundamental rights

4 Ethical Principles

Acknowledge and address tensions between them

- Respect for human autonomy
- Prevention of harm
- Fairness
- Explicability

CHAPTER I

Realisation of Trustworthy AI

Implement the key requirements

7 Key Requirements

Evaluate and address these continuously throughout the AI system's life cycle

via

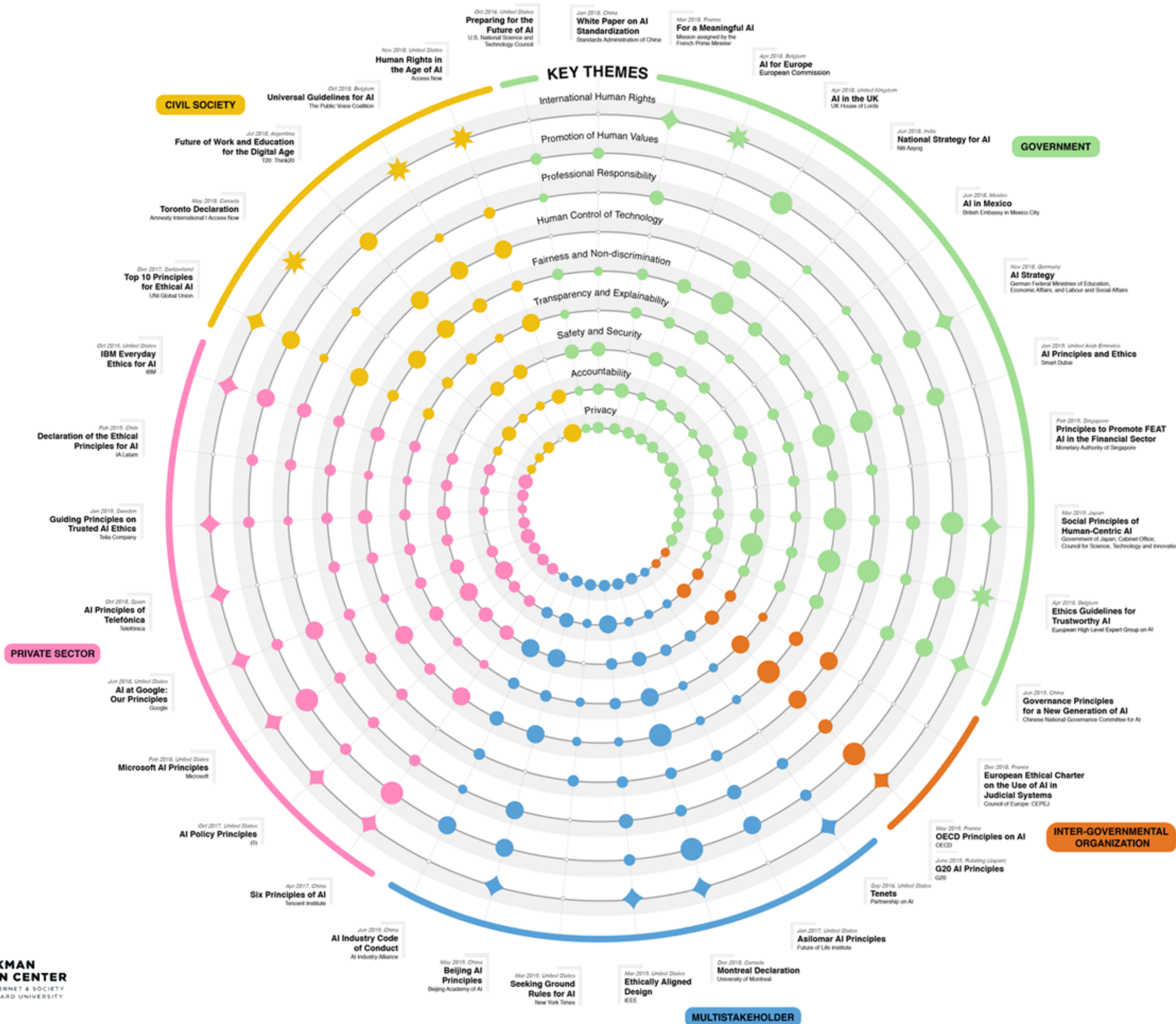
**Technical
Methods**

**Non-Technical
Methods**

- Human agency and oversight
- Technical robustness and safety
- Privacy and data governance
- Transparency
- Diversity, non-discrimination and fairness
- Societal and environmental wellbeing
- Accountability

CHAPTER II

Proliferation of AI ethics principles...



Consolidated				Ethics Guidelines for Trustworthy AI (EC)	Survey of Ethics Principles (Jobin et al 2019)
Individual			AUTONOMY	HUMAN AGENCY & OVERSIGHT	FREEDOM & AUTONOMY
			DIGNITY		DIGNITY
			PRIVACY	PRIVACY & DATA GOVERNANCE	PRIVACY
			FAIRNESS	DIVERSITY, NON-DISCRIMINATION, FAIRNESS	JUSTICE & FAIRNESS
Social			SOLIDARITY		SOLIDARITY
			SOCIAL WELLBEING	SOCIETAL & ENVIRONMENTAL WELLBEING	SUSTAINABILITY
					BENEFICENCE & NON-MALEFICENCE
			PERFORMANCE		TRUST
Technical			SAFETY	TECHNICAL ROBUSTNESS & SAFETY	
			EXPLAINABILITY & ACCOUNTABILITY	TRANSPARENCY	TRANSPARENCY
				ACCOUNTABILITY	RESPONSIBILITY

Dilemmas in AI ethics: True, Practical

True: Privacy versus Fairness

Perfectly safeguarding users' personal information prohibits assessments of AI functioning across distinct demographic groups.

Practical: Accuracy versus explainability

The most accurate AI results can emerge from an internal logic that developers or users do not fully understand. Still, perfect accuracy can coexist with full explainability.