

**Meeting of the Council at Ministerial Level, 2-3 May 2024****REVISED RECOMMENDATION OF THE COUNCIL ON ARTIFICIAL  
INTELLIGENCE****(Adopted by the Council at Ministerial level on 3 May 2024)****JT03542983**

**THE COUNCIL,**

**HAVING REGARD** to Article 5 b) of the Convention on the Organisation for Economic Co-operation and Development of 14 December 1960;

**HAVING REGARD** to standards developed by the OECD in the areas of privacy, digital security, consumer protection and responsible business conduct;

**HAVING REGARD** to the Sustainable Development Goals set out in the 2030 Agenda for Sustainable Development adopted by the United Nations General Assembly (A/RES/70/1) as well as the 1948 Universal Declaration of Human Rights;

**HAVING REGARD** to the important work being carried out on artificial intelligence (hereafter, “AI”) in other international governmental and non-governmental fora;

**RECOGNISING** that AI has pervasive, far-reaching and global implications that are transforming societies, economic sectors and the world of work, and are likely to increasingly do so in the future;

**RECOGNISING** that AI has the potential to improve the welfare and well-being of people, to contribute to positive sustainable global economic activity, to increase innovation and productivity, and to help respond to key global challenges;

**RECOGNISING** that, at the same time, these transformations may have disparate effects within, and between societies and economies, notably regarding economic shifts, competition, transitions in the labour market, inequalities, and implications for democracy and human rights, privacy and data protection, and digital security;

**RECOGNISING** that trust is a key enabler of digital transformation; that, although the nature of future AI applications and their implications may be hard to foresee, the trustworthiness of AI systems is a key factor for the diffusion and adoption of AI; and that a well-informed whole-of-society public debate is necessary for capturing the beneficial potential of the technology, while limiting the risks associated with it;

**UNDERLINING** that certain existing national and international legal, regulatory and policy frameworks already have relevance to AI, including those related to human rights, consumer and personal data protection, intellectual property rights, responsible business conduct, and competition, while noting that the appropriateness of some frameworks may need to be assessed and new approaches developed;

**RECOGNISING** that given the rapid development and implementation of AI, there is a need for a stable policy environment that promotes a human-centric approach to trustworthy AI, that fosters research, preserves economic incentives to innovate, and that applies to all stakeholders according to their role and the context;

**CONSIDERING** that embracing the opportunities offered, and addressing the challenges raised, by AI applications, and empowering stakeholders to engage is essential to fostering adoption of trustworthy AI in society, and to turning AI trustworthiness into a competitive parameter in the global marketplace.

**On the proposal of the Digital Policy Committee:**

**I. AGREES** that for the purpose of this Recommendation the following terms should be understood as follows:

- *AI system*: An AI system is a machine-based system that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments. Different AI systems vary in their levels of autonomy and

adaptiveness after deployment.

- *AI system lifecycle*: An AI system lifecycle typically involves several phases that include to: plan and design; collect and process data; build model(s) and/or adapt existing model(s) to specific tasks; test, evaluate, verify and validate; make available for use/deploy; operate and monitor; and retire/decommission. These phases often take place in an iterative manner and are not necessarily sequential. The decision to retire an AI system from operation may occur at any point during the operation and monitoring phase.
- *AI actors*: AI actors are those who play an active role in the AI system lifecycle, including organisations and individuals that deploy or operate AI.
- *AI knowledge*: AI knowledge refers to the skills and resources, such as data, code, algorithms, models, research, know-how, training programmes, governance, processes, and best practices required to understand and participate in the AI system lifecycle, including managing risks.
- *Stakeholders*: Stakeholders encompass all organisations and individuals involved in, or affected by, AI systems, directly or indirectly. AI actors are a subset of stakeholders.

### **Section 1: Principles for responsible stewardship of trustworthy AI**

**II. RECOMMENDS** that Members and non-Members adhering to this Recommendation (hereafter the “Adherents”) promote and implement the following principles for responsible stewardship of trustworthy AI, which are relevant to all stakeholders.

**III. CALLS ON** all AI actors to promote and implement, according to their respective roles, the following principles for responsible stewardship of trustworthy AI.

**IV. UNDERLINES** that the following principles are complementary and should be considered as a whole.

#### **1.1. Inclusive growth, sustainable development and well-being**

Stakeholders should proactively engage in responsible stewardship of trustworthy AI in pursuit of beneficial outcomes for people and the planet, such as augmenting human capabilities and enhancing creativity, advancing inclusion of underrepresented populations, reducing economic, social, gender and other inequalities, and protecting natural environments, thus invigorating inclusive growth, well-being, sustainable development and environmental sustainability.

#### **1.2. Respect for the rule of law, human rights and democratic values, including fairness and privacy**

- a) AI actors should respect the rule of law, human rights, democratic and human-centred values throughout the AI system lifecycle. These include non-discrimination and equality, freedom, dignity, autonomy of individuals, privacy and data protection, diversity, fairness, social justice, and internationally recognised labour rights. This also includes addressing misinformation and disinformation amplified by AI, while respecting freedom of expression and other rights and freedoms protected by applicable international law.
- b) To this end, AI actors should implement mechanisms and safeguards, such as capacity for human agency and oversight, including to address risks arising from uses outside of intended purpose, intentional misuse, or unintentional misuse in a manner appropriate to the context and consistent with the state of the art.

### **1.3. Transparency and explainability**

AI Actors should commit to transparency and responsible disclosure regarding AI systems. To this end, they should provide meaningful information, appropriate to the context, and consistent with the state of art:

- i. to foster a general understanding of AI systems, including their capabilities and limitations,
- ii. to make stakeholders aware of their interactions with AI systems, including in the workplace,
- iii. where feasible and useful, to provide plain and easy-to-understand information on the sources of data/input, factors, processes and/or logic that led to the prediction, content, recommendation or decision, to enable those affected by an AI system to understand the output, and,
- iv. to provide information that enable those adversely affected by an AI system to challenge its output.

### **1.4. Robustness, security and safety**

- a) AI systems should be robust, secure and safe throughout their entire lifecycle so that, in conditions of normal use, foreseeable use or misuse, or other adverse conditions, they function appropriately and do not pose unreasonable safety and/or security risks.
- b) Mechanisms should be in place, as appropriate, to ensure that if AI systems risk causing undue harm or exhibit undesired behaviour, they can be overridden, repaired, and/or decommissioned safely as needed.
- c) Mechanisms should also, where technically feasible, be in place to bolster information integrity while ensuring respect for freedom of expression.

### **1.5. Accountability**

- a) AI actors should be accountable for the proper functioning of AI systems and for the respect of the above principles, based on their roles, the context, and consistent with the state of the art.
- b) To this end, AI actors should ensure traceability, including in relation to datasets, processes and decisions made during the AI system lifecycle, to enable analysis of the AI system's outputs and responses to inquiry, appropriate to the context and consistent with the state of the art.
- c) AI actors, should, based on their roles, the context, and their ability to act, apply a systematic risk management approach to each phase of the AI system lifecycle on an ongoing basis and adopt responsible business conduct to address risks related to AI systems, including, as appropriate, via co-operation between different AI actors, suppliers of AI knowledge and AI resources, AI system users, and other stakeholders. Risks include those related to harmful bias, human rights including safety, security, and privacy, as well as labour and intellectual property rights.

## **Section 2: National policies and international co-operation for trustworthy AI**

**V. RECOMMENDS** that Adherents implement the following recommendations, consistent with the principles in section 1, in their national policies and international co-operation, with special attention to small and medium-sized enterprises (SMEs).

## 2.1. Investing in AI research and development

- a) Governments should consider long-term public investment, and encourage private investment, in research and development and open science, including interdisciplinary efforts, to spur innovation in trustworthy AI that focus on challenging technical issues and on AI-related social, legal and ethical implications and policy issues.
- b) Governments should also consider public investment and encourage private investment in open-source tools and open datasets that are representative and respect privacy and data protection to support an environment for AI research and development that is free of harmful bias and to improve interoperability and use of standards.

## 2.2. Fostering an inclusive AI-enabling ecosystem

Governments should foster the development of, and access to, an inclusive, dynamic, sustainable, and interoperable digital ecosystem for trustworthy AI. Such an ecosystem includes *inter alia*, data, AI technologies, computational and connectivity infrastructure, and mechanisms for sharing AI knowledge, as appropriate. In this regard, governments should consider promoting mechanisms, such as data trusts, to support the safe, fair, legal and ethical sharing of data.

## 2.3. Shaping an enabling interoperable governance and policy environment for AI

- a) Governments should promote an agile policy environment that supports transitioning from the research and development stage to the deployment and operation stage for trustworthy AI systems. To this effect, they should consider using experimentation to provide a controlled environment in which AI systems can be tested, and scaled-up, as appropriate. They should also adopt outcome-based approaches that provide flexibility in achieving governance objectives and co-operate within and across jurisdictions to promote interoperable governance and policy environments, as appropriate.
- b) Governments should review and adapt, as appropriate, their policy and regulatory frameworks and assessment mechanisms as they apply to AI systems to encourage innovation and competition for trustworthy AI.

## 2.4. Building human capacity and preparing for labour market transformation

- a) Governments should work closely with stakeholders to prepare for the transformation of the world of work and of society. They should empower people to effectively use and interact with AI systems across the breadth of applications, including by equipping them with the necessary skills.
- b) Governments should take steps, including through social dialogue, to ensure a fair transition for workers as AI is deployed, such as through training programmes along the working life, support for those affected by displacement, including through social protection, and access to new opportunities in the labour market.
- c) Governments should also work closely with stakeholders to promote the responsible use of AI at work, to enhance the safety of workers, the quality of jobs and of public services, to foster entrepreneurship and productivity, and aim to ensure that the benefits from AI are broadly and fairly shared.

## 2.5. International co-operation for trustworthy AI

- a) Governments, including developing countries and with stakeholders, should actively co-operate to advance these principles and to progress on responsible stewardship of trustworthy AI.

- b) Governments should work together in the OECD and other global and regional fora to foster the sharing of AI knowledge, as appropriate. They should encourage international, cross-sectoral and open multi-stakeholder initiatives to garner long-term expertise on AI.
- c) Governments should promote the development of multi-stakeholder, consensus-driven global technical standards for interoperable and trustworthy AI.
- d) Governments should also encourage the development, and their own use, of internationally comparable indicators to measure AI research, development and deployment, and gather the evidence base to assess progress in the implementation of these principles.

**VI. INVITES** the Secretary-General and Adherents to disseminate this Recommendation.

**VII. INVITES** non-Adherents to take due account of, and adhere to, this Recommendation.

**VIII. INSTRUCTS** the Digital Policy Committee, through its Working Party on AI Governance, to:

- a) continue its important work on artificial intelligence building on this Recommendation and taking into account work in other international fora, and to further develop the measurement framework for evidence-based AI policies;
- b) develop and iterate further practical guidance on the implementation of this Recommendation to meet evolving developments and new policy priorities;
- c) provide a forum for exchanging information on AI policy and activities including experience with the implementation of this Recommendation, and to foster multi-stakeholder and interdisciplinary dialogue to promote trust in and adoption of AI; and
- d) report to Council, in consultation with other relevant committees, on the implementation, dissemination and continued relevance of this Recommendation no later than five years following its revision and at least every ten years thereafter.