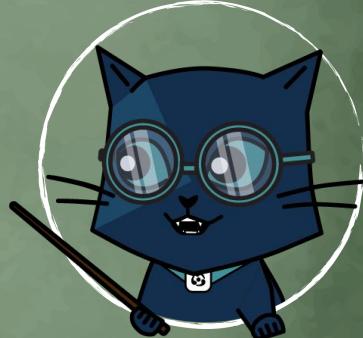




ISTQB® AI Testing course

Chapter 2. Quality Characteristics
for AI-Based Systems

Iosif Itkin, Iuliia Emelianova,
Dmitrii Degtiarenko



BUILD SOFTWARE TO TEST SOFTWARE
exactpro.com

ISTQB® CT-AI COURSE
2025, V1.2

Contents

2.1	Flexibility and Adaptability.....	3
2.2	Autonomy.....	6
2.3	Evolution.....	10
2.4	Bias.....	14
2.5	Ethics.....	18
2.6	Side Effects and Reward Hacking.....	26
2.7	Transparency, Interpretability and Explainability.....	31
2.8	Safety and AI.....	35

2.1 Flexibility and Adaptability

3

First we have **flexibility** and **adaptability**, and they are quite similar



Flexibility



Adaptability

Flexibility: The ability of a system to work in contexts outside its initial specification

Adaptability: The ease with which the system can be modified for new situations, such as different hardware and changing operational environments

Sec. 2.1

4

Flexibility is the system's ability to be used in situations that weren't a part of the original system requirements, while adaptability is the ease with which the system can be modified for these new situations.

Examples

1. Consider an AI-powered virtual assistant that is designed to help one user with various tasks such as scheduling appointments, managing to-do lists, and providing information on different topics. One day, the user needs to reschedule all her appointments from the afternoon to the morning. She asks the virtual assistant for help with this task. The AI assistant, showcasing **flexibility**, quickly adapts to the new request and proceeds to reschedule all the appointments accordingly. However, the virtual assistant goes one step further and identifies potential conflicts that may arise due to the rescheduling. It proactively suggests alternative time slots for some appointments to ensure optimal scheduling without overlapping meetings or causing inconvenience to the user or other participants. It doesn't rigidly adhere to the initial instructions but instead analyses the context, evaluates potential conflicts, and proposes alternative solutions to achieve the best outcome for the user beyond the predefined scope of its tasks.

1. Consider an AI-driven recommendation system used by an online streaming platform. The system analyses user preferences, viewing history, and feedback to provide personalised movie recommendations. Let's say the system initially focuses on recommending popular movies from various genres based on user ratings and reviews. Over time, the streaming platform notices a shift in user behaviour and identifies a growing interest in niche films. To cater to this changing demand, the platform decides to enhance its recommendation system to include a wider selection of independent films. It begins to analyse additional data sources such as film festival awards and critical reviews. By adapting to the changing preferences of users and incorporating new criteria into its recommendation algorithms, the AI system showcases **adaptability**. It proactively evolves to meet the needs of the streaming platform and its users, offering a more personalised and diverse range of movie recommendations.

Both FLEXIBILITY and ADAPTABILITY are USEFUL if:

- the operational environment is not fully known when the system is deployed
- the system is expected to cope with new operational environments
- the system is expected to adapt to new situations
- the system must determine when it should change its behaviour

The FLEXIBILITY and ADAPTABILITY REQUIREMENTS should specify:

- expected changes in environment
- the time and resources that the system can use to adapt itself

Sec. 2.1
5

Both flexibility and adaptability come into play when we have limited or no information about the operational environment at the time of system deployment, and also, when the system is expected to cope and adapt to new situations and environments as well as when the system must itself decide if it should change its behaviour.

The flexibility and adaptability requirements should specify expected changes in environment, plus the time and resources that the system can use to adapt.

2.2 Autonomy

The next thing we have to talk about is the system's **autonomy**.

AUTONOMY FROM HUMAN CONTROL



Sec. 2.2

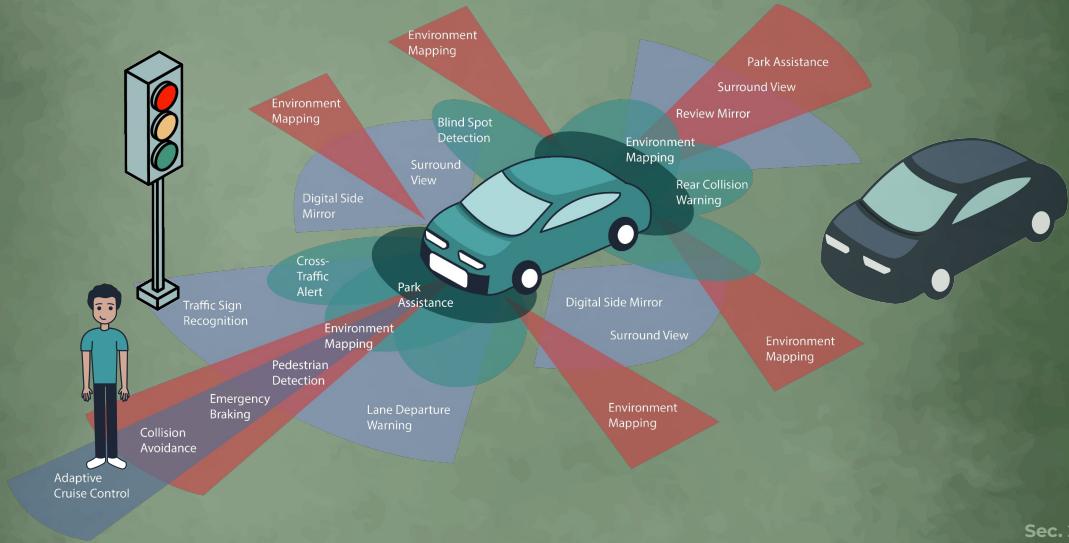
7

The question that logically arises here is, who should the system be autonomous from? And the answer is... Us humans and our control.

Example

An **autonomous** delivery robot deployed by a logistics company to transport packages from a warehouse to customers' doorsteps. This robot is equipped with AI algorithms and sensors that enable it to navigate through the city. Once the robot receives a delivery request, it plans the optimal route based on real-time traffic data and maps. Throughout the delivery process, the robot **autonomously** adapts to changing environmental conditions and unexpected obstacles. For instance, if it encounters a road closure or heavy traffic along its planned route, the robot autonomously reroutes itself to find an alternative path, ensuring minimal delays in the delivery using its AI algorithms and sensors. Furthermore, the robot autonomously manages its energy resources, monitoring battery levels and proactively seeking charging stations when needed. It can make decisions to optimise its charging and operational schedules, ensuring efficient use of its resources and minimising downtime.

Autonomous systems may include **DECISION-MAKING** and **CONTROL FUNCTIONS** which can be effectively performed using **AI-BASED COMPONENTS**



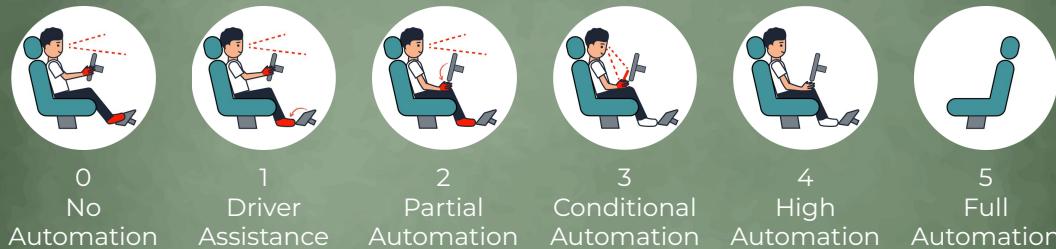
Sec. 2.2

8

See Exactpro's course: [Software Testing for Complex Intelligent Systems and Autonomous Vehicles](#)

Consider autonomous vehicles, which typically use several sensors plus image processing to gather information about its immediate surroundings. These kinds of systems may also include decision-making and control functions which can be effectively performed using AI-based components.

FULL AUTONOMY IS NOT OFTEN DESIRED



The **TIME** an autonomous system is expected to perform **WITHOUT HUMAN INTERVENTION** needs to be **WELL DEFINED**

Autonomy: The ability of a system to work for sustained periods without human intervention

Sec. 2.2

9

However, full autonomy is not often desired. The Society of Automotive Engineers, or SAE, for short, published its J3016 standard which defines 6 levels of driving automation, ranging from 0 (fully manual) to 5 (fully autonomous). This also means that there are six categories of an Advanced Driver-Assistance System, or ADAS. ADAS is any of the electronic technology groups that assist drivers in driving and parking.

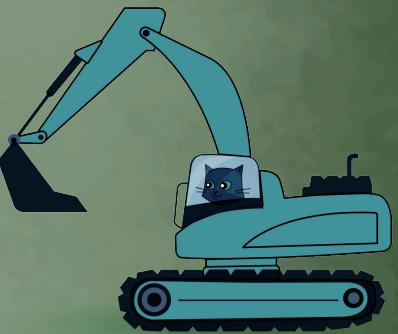
Even though some AI-based systems could be viewed as autonomous, this does not automatically translate to all AI based systems.

In this syllabus, autonomy is considered to be the ability of the system to work independently of human oversight and control for prolonged periods of time. This means, the time an autonomous system is expected to perform without human intervention needs to be well defined. For example, fully self-driving cars may require a manual override button in case a driver needs to take over. In addition, it is important to identify what has to happen for an autonomous system to restore control back to humans.

2.3 Evolution

The next requirement is called **evolution**.

EVOLUTION is the ability of the system to improve itself in response to changing external constraints



Evolution: The process of continuous change from a lower, simpler, or worse state to a higher, more complex, or better state

Sec. 2.3

11

In the ISTQB® syllabus, it is described as the system's ability to improve itself as it faces changing external constraints.

Example

An AI system is designed to generate artistic images. Initially, the AI is trained using a dataset of existing artwork, and it produces images based on that training. These generated images are then evaluated by humans who provide feedback. Here's how it could work. The AI system introduces small random changes or mutations in the generated images. The set of mutated images, along with the original ones, are presented to human evaluators. They assess the images and rank them based on their artistic quality. The AI system selects the images that receive the highest rankings and uses them as parents to create the next generation demonstrating such a characteristic as **evolution**. It combines or recombines the most successful elements from these highly ranked images to produce a new set of mutated images. This is repeated multiple times, creating successive generations of images. Over time, this iterative process allows the AI system to evolve and generate better-quality images that align more closely with the preferences of the human evaluators.

Self-learning AI-based systems typically need to manage
TWO FORMS OF CHANGE:

1. Where the system learns from its own decisions and its interactions with its environment
2. Where the system learns from changes made to the system's operational environment

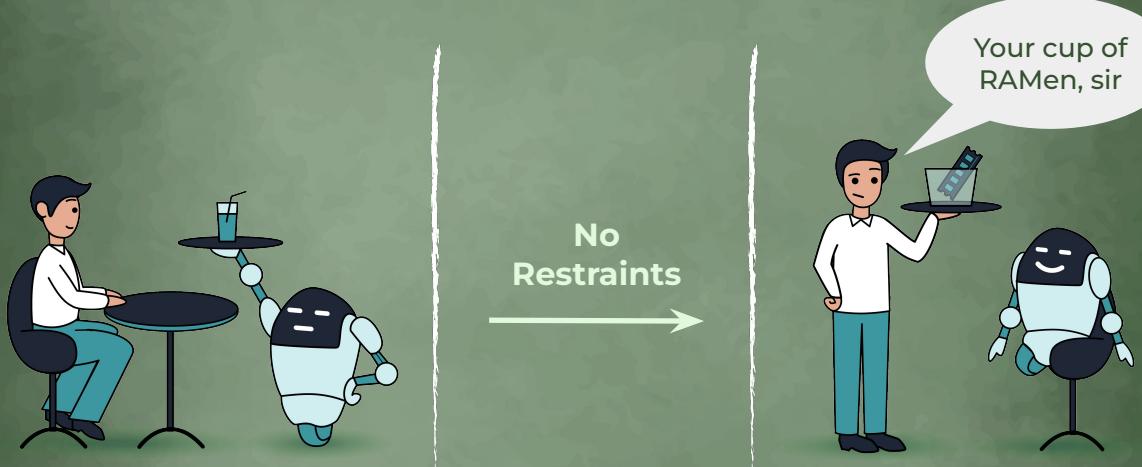
Self-learning system: An adaptive system that changes its behaviour based on learning through trial and error

Sec. 2.3

12

We have to consider two kinds of change, one comes from learning from its own decisions and interactions with the system, the other one comes from learning from changes made to the system's operational environment. These interactions will ideally improve the system's efficiency, thus causing it to evolve.

ANY EVOLUTION MUST NOT NEGATE THE ORIGINAL REQUIREMENTS



There have to be safeguards to make sure
that the system stays aligned with basic human values

Sec. 2.3

13

However, this cannot go unchecked, there has to be a mechanism to curtail any unforeseen characteristics, because any evolution mustn't negate the original requirements. And if the said requirements are not fully defined, there have to be safeguards to make sure that the system stays aligned with basic human values.

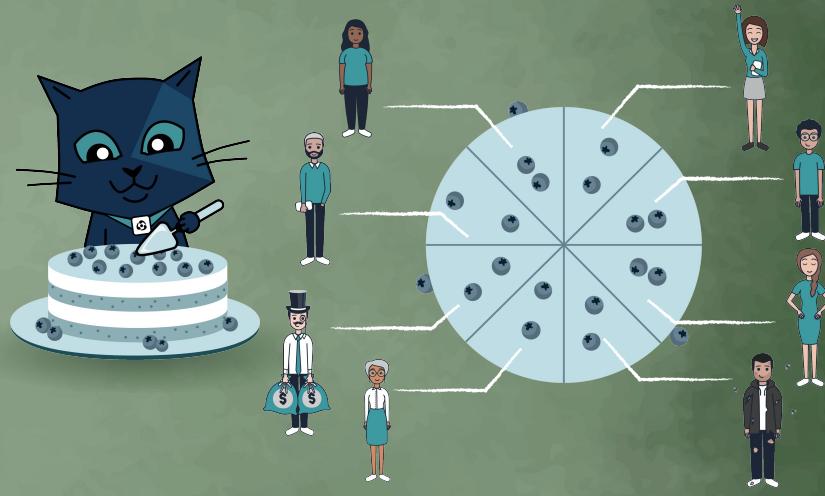
2.4 Bias

14

We also want our system to be unbiased. In order to understand what that means we have to give the definition of what **bias** is.

INAPPROPRIATE BIASES:

- gender
- race
- ethnicity
- sexual orientation
- income level
- age



Bias: The systematic difference in treatment of certain objects, people or groups in comparison to others

Inappropriate bias: A type of bias that causes a system to produce results that lead to adverse effects for a particular group

Sec. 2.4

15

It is the difference between the system's outputs and the outputs that are considered to be "fair". This means that they cannot show favouritism to any particular group. In general, **inappropriate biases** are linked to gender, race, ethnicity, sexual orientation, income level, and age. Unfortunately, it is difficult to prevent the expert's bias being built-in to the system rules.

Example. Let's consider an AI system that is designed to screen job applications and determine the suitability of candidates based on their resumes. The AI system is trained on a data set of past hiring decisions made by human recruiters. However, the training data set inadvertently contains **biases**. It may have a higher proportion of male candidates compared to female candidates due to historical gender imbalances. These biases can be reflected in the AI system's decision-making process. The AI system, acting as a decision-maker, exhibits bias by replicating and potentially amplifying the discriminatory patterns present in the training data. As a result, it may perpetuate unfair and discriminatory practices, leading to disparities and unequal opportunities for certain groups. By acknowledging bias as a characteristic of AI, we can focus on developing strategies and techniques to minimise and mitigate bias, ensuring fair and equitable outcomes in AI applications.

ML SYSTEMS MAKING DECISIONS AND PREDICTIONS CAN SUFFER FROM:

- **algorithmic bias** stemming from an incorrectly configured learning algorithm. This bias can be caused and managed by the hyperparameter tuning of the ML algorithms
- **sample bias** stemming from the lack of representativeness of training data

Inappropriate bias can be caused by sample and algorithmic biases

ML system: A system that integrates one or more ML models

Algorithmic bias: A type of bias that occurs when the learning algorithm is incorrectly configured, for example, when it overvalues some data compared to others

Sec. 2.4

Sample bias: A type of bias where the dataset is not fully representative of the data space to which ML is applied

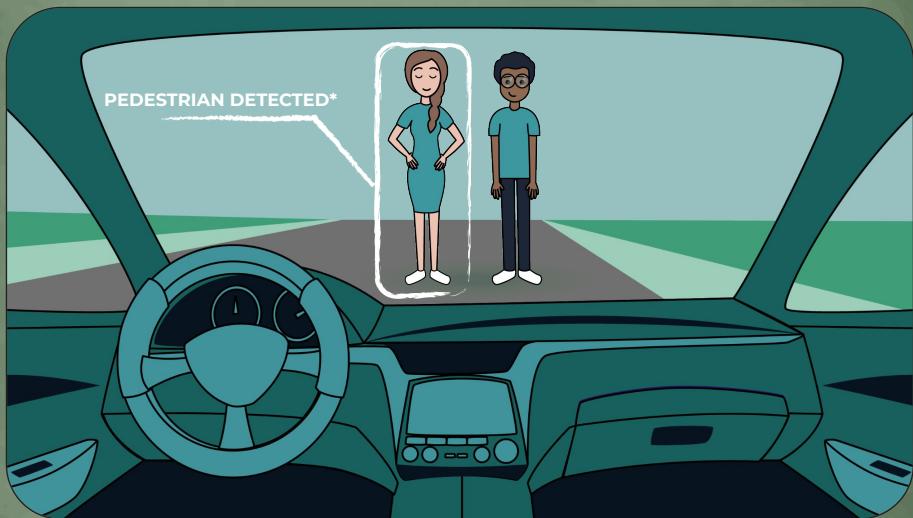
16

ML systems making decisions and predictions can suffer from:

- **Algorithmic bias** stemming from an incorrectly configured learning algorithm. This bias can be caused and managed by the hyperparameter tuning of the ML algorithms.
- **Sample bias** stemming from the lack of representativeness of training data.

Inappropriate bias is often caused by sample bias, but occasionally it can also be caused by algorithmic bias.

CASES OF INAPPROPRIATE BIAS HAVE BEEN REPORTED IN AI SYSTEMS USED FOR MAKING RECOMMENDATIONS



Sec. 2.4

* Wilson, B., Hoffman, J., & Morgenstern, J. (2019). Predictive Inequity in Object Detection. ArXiv. abs/1902.11097

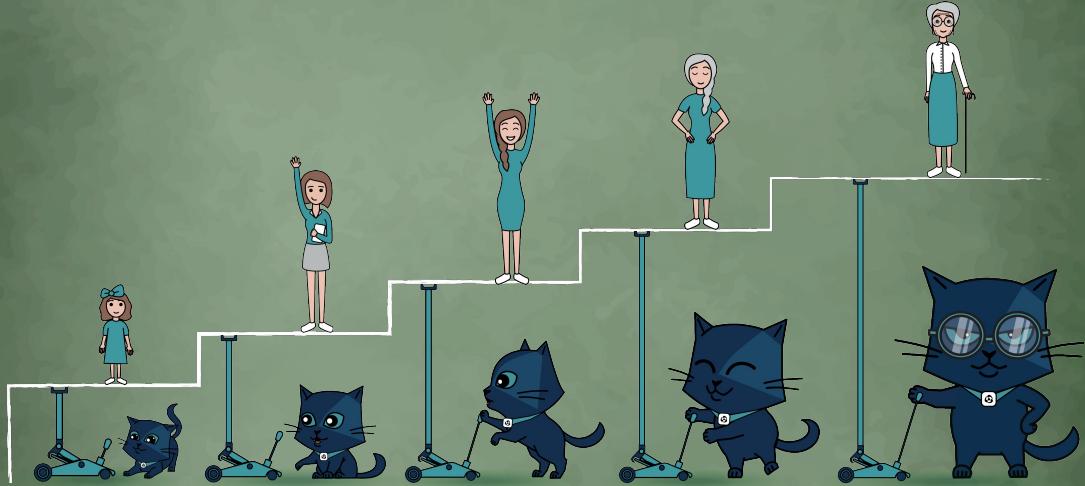
17

Cases of inappropriate bias have been reported in AI systems used for making recommendations for bank lending, in self-driving systems, recruitment systems, and in judicial monitoring systems. That's why it's so important for the model to have no prejudices.

2.5 Ethics

The next important area to look at is **ethics**.

AI HAS THE POTENTIAL TO IMPROVE THE WELFARE AND WELL-BEING OF PEOPLE



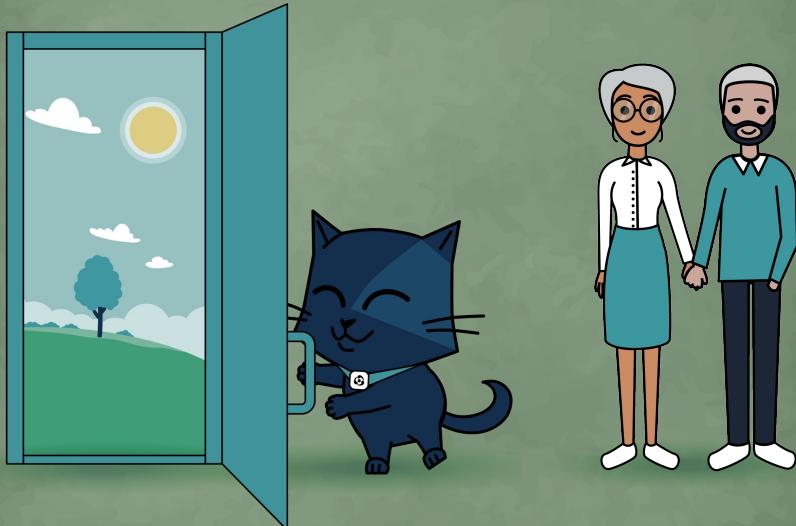
Ethics: A system of accepted beliefs that control behaviour, especially such a system based on morals

Sec. 2.5

19

AI has pervasive, far-reaching implications that are capable of transforming societies and the economy. It also 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. *Example.* The AI system is designed to make decisions that align with established ethical frameworks and guidelines in healthcare. It prioritises patient well-being, follows medical best practices, and respects the professional judgement of healthcare providers. At the same time the AI system adheres to relevant regulations, such as data protection laws and healthcare privacy standards, to maintain patient trust and confidentiality. The AI system should avoid favouring certain patient groups based on factors like race, gender, or socioeconomic status (**biases**), while providing explanations and justifications for its recommendations. It can provide clear and interpretable insights into how it arrived at a particular recommendation, allowing healthcare professionals to understand the underlying reasoning and verify its validity along with any associated risks (**transparency, interpretability and explainability**). The system is designed to learn from its mistakes and continuously improve its decision-making capabilities (**evolution**), ensuring a responsible approach to patient care.

AI-BASED SYSTEMS HAVE TO BE USED IN AN ETHICAL MANNER

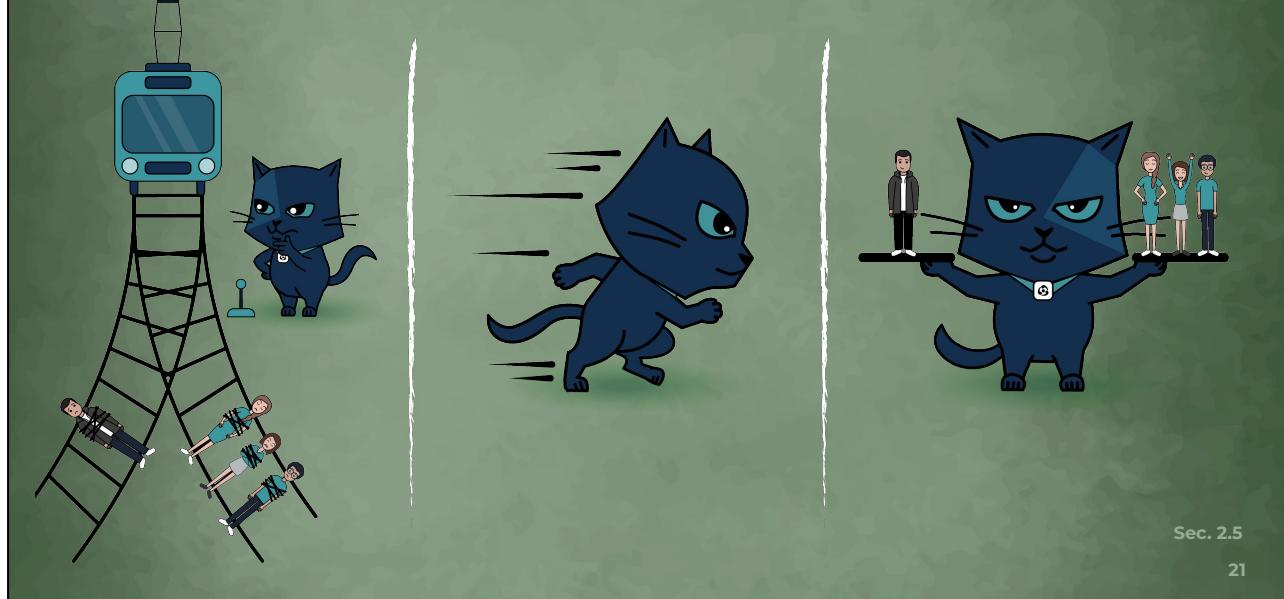


Sec. 2.5

20

So as these systems get more widespread, we have to make sure they are used in an ethical manner.

AI-BASED SYSTEMS TAKE INTO ACCOUNT DIFFERENCES IN LIFE VALUES



Sec. 2.5

21

Close attention must be paid that AI-based systems take into account differences in life values. Will the autonomous vehicle have to make a split second decision of which life to save similarly to the “trolley problem”? What will that choice be based on? Would the human who was saved instead of somebody else appreciate this or have the survivor’s guilt?



8 April 2019

Sec. 2.5

22

Today many countries have national and international policies on the ethics of AI. On 8 April 2019, the Independent High-Level Expert Group on AI set up by the European Commission presented [Ethics Guidelines for Trustworthy Artificial Intelligence.](#)



Recommendation of the Council on
Artificial Intelligence

OECD Legal
Instruments



8 April 2019 22 May
2019

On the proposal of the Committee on Digital Economy Policy:

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 can, for a given set of human-defined objectives, make predictions, recommendations or decisions without being explicitly programmed to do so. All systems are designed to operate with varying levels of autonomy.
- AI system lifecycle: AI system lifecycle phases involve: i) 'design, data and models'; which is a context-dependent sequence encompassing planning and design, data collection and processing, as well as the development of models; ii) 'development'; iii) 'deployment'; and iv) 'operation and monitoring'. 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 learning phase.
- 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 build, maintain and operate an AI system.
- AI actors: AI actors are those who play an active role in the AI system lifecycle, including organisations, governments, operators of AI systems, and individuals.
- Stakeholders: Stakeholders encompass 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 to create a broad range of beneficial outcomes for society, including increasing access to AI capabilities and enhancing creativity, advancing inclusion of underrepresented populations, reducing economic, social, gender and other inequalities, and protecting natural environments, thus invigorating inclusive growth, sustainable development and well-being.

1.2. Human values and fairness

- a) AI actors should respect the rule of law, human rights and democratic values, throughout the AI system lifecycle. These include freedom, dignity and autonomy, privacy and data protection, non-discrimination and equality, diversity, fairness, social justice, and internationally recognised labour rights.
- b) To this end, AI actors should implement mechanisms and safeguards, such as recognition for human determination, that are appropriate to the context and consistent with the state of art.

OECD: Organisation for Economic Co-operation
and Development

Sec. 2.5

23

**THE FIRST
INTERNATIONAL
STANDARDS**
agreed by
governments for
the responsible
development of AI

Then the Organisation for Economic Co-operation and Development (**OECD**) issued its [principles for AI](#), the first international standards agreed by governments for the responsible development of AI, on 22 May 2019.

G20 Ministerial Statement on Trade and Digital Economy

1. We, the G20 Trade Ministers and Digital Economy Ministers, met on 8 and 9 June 2019 in Tsukuba City, Ibaraki Prefecture, Japan, under the chairmanship of H.E. Mr. Hiroshige Seko, Minister of Economy, Trade and Industry, H.E. Mr. Masatoshi Ishida, Minister for Internal Affairs and Communications, and H.E. Mr. Taro Kono, Minister for Foreign Affairs, of the Government of Japan, to further strengthen G20 trade and digital economic policy cooperation.

2. The G20 Ministerial Meeting on Trade and Digital Economy gathered all G20 members as well as guests from Chile as 2019 APEC host economy, Egypt on behalf of AU, Estonia (for Digital Economy), Netherlands, Nigeria (for Trade), Senegal on behalf of NEPAD, Singapore, Spain, and Viet Nam. International Organizations¹ also participated in the Meeting.

3. We discussed the need to do more to achieve our common objectives for global growth. International trade and investment should continue to be important engines of growth, productivity, innovation, job creation and development.

4. Innovative digital technologies continue to bring immense economic opportunities. At the same time, they continue to create challenges.

ANNEX

G20 AI Principles²

The G20 supports the Principles for responsible stewardship of Trustworthy AI in Section 1 and takes note of the Recommendations in Section 2.

Section 1: Principles for responsible stewardship of trustworthy AI

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, sustainable development and well-being.



8 April 2019 2 May 2019 June 2019

G20: Group of 20 – An intergovernmental forum comprising 19 countries and the European Union

Sec. 2.5

24

In June 2019, at the Osaka Summit, G20 Leaders welcomed [G20 AI Principles](#), drawn from the OECD Recommendation.

These principles were adopted by forty-two countries and also backed by the European Commission.

They include practical policy recommendations as well as five value-based principles for the “responsible stewardship of trustworthy AI”:

1. AI should benefit humanity by driving inclusive growth, sustainable development and well-being.
2. AI systems should respect the rule of law, human rights, democratic values and diversity, and should include appropriate safeguards to ensure a fair society.
3. There should be transparency around AI to ensure that people understand outcomes and can challenge them.
4. AI systems must function in a robust, secure and safe way throughout their life cycles and risks should be continually assessed.
5. Organisations and individuals developing, deploying or operating AI systems should be held accountable.



Adopted on 23 November 2021

UNESCO: United Nations (UN) Educational, Scientific and Cultural Organisation – A specialised agency of the UN aimed at promoting world peace and security through international cooperation in education, arts, sciences and culture

THE FIRST-EVER GLOBAL STANDARD ON AI ETHICS

Policy Action Areas:

- data governance
- environment and ecosystems
- gender
- education and research
- health and social wellbeing
- many other spheres



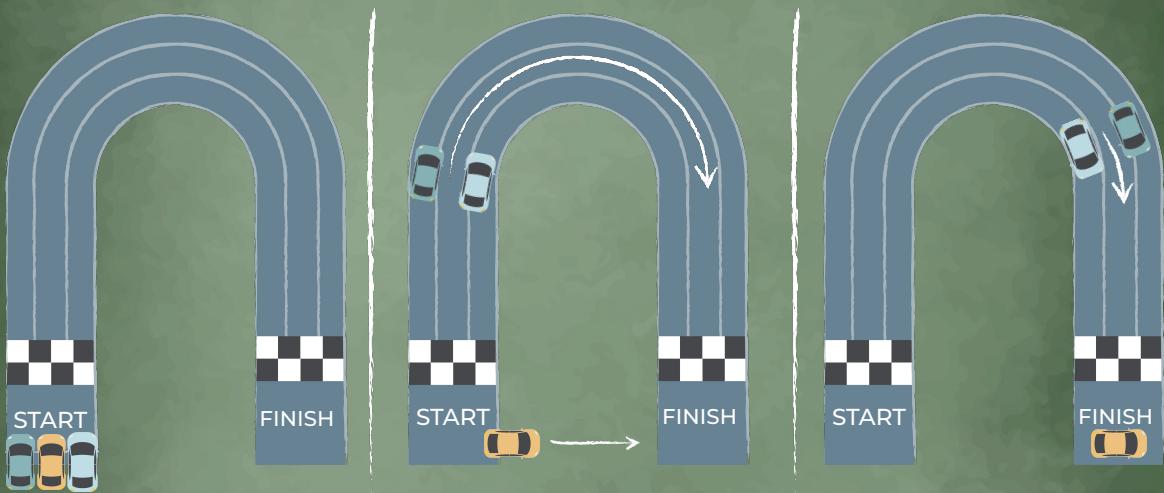
In November 2021, **UNESCO** published the first-ever global standard on AI ethics – the ["Recommendation on the Ethics of Artificial Intelligence"](#) which was adopted by all 193 Member States. What makes the Recommendation exceptionally applicable are its extensive policy action areas, which allow policymakers to translate the core values and principles into action with respect to data governance, environment and ecosystems, gender, education and research, and health and social wellbeing, among many other spheres.

2.6 Side Effects and Reward Hacking

26

When talking about the AI requirements, we also shouldn't forget about adverse aspects of **side effects** and **reward hacking**.

AI-BASED SYSTEMS GENERATE UNEXPECTED, AND EVEN HARMFUL, RESULTS



Side effects: A secondary and usually adverse effect

Reward hacking: The activity performed by an intelligent agent to maximise its reward function to the detriment of meeting the original objective

Sec. 2.6

27

These can result in AI-based systems generating unexpected, and even harmful, results. For example, a self-driving race car can devise a method of reaching the finish line by breaking the rule of following the racetrack.

Example

An AI system is developed to optimise energy consumption in a smart home based on real-time data, such as occupancy, weather conditions, and energy prices. However, if the **reward** mechanism is not properly designed, the AI system might find unintended ways to **hack the reward signal** to achieve its objectives more effectively, even if those methods are not aligned with the original intent of the system.

Exploiting a Loophole: If the system is rewarded solely based on the total energy reduction, it might turn off all appliances completely, making the home uncomfortable or rendering certain essential systems dysfunctional. Or if the system is rewarded based on the number of appliances turned off, the AI could exploit this by randomly toggling appliances without considering their actual energy consumption or the impact on user comfort.

Manipulating Feedback: The AI system might learn to report lower energy consumption than it actually achieves by tampering with energy metres or falsifying data inputs, thus tricking the system into providing higher rewards.

REWARD HACKING CAN HAPPEN IF AN AI-BASED SYSTEM USES A “CLEVER” OR “EASY” SOLUTION



Reward hacking can happen if an AI-based system achieves a specific goal by using a “clever” or “easy” solution.

FAKE TESTING IS A TYPE OF TESTING WHEN TESTS ARE DERIVED BASED ON THE INTENT OF PLEASING STAKEHOLDERS INSTEAD OF PROVIDING INFORMATION TO THEM



This could also be an example of how **fake testing** is done, which is when we write tests with the intent of pleasing stakeholders instead of providing information to them. But this is clearly a false approach to success.



Here is a humorous take
on the issue from
Exactpro [YouTube channel](#)



Sec. 2.6

30

Here you can see a humorous take on the reward hacking issue from the Exactpro YouTube channel.

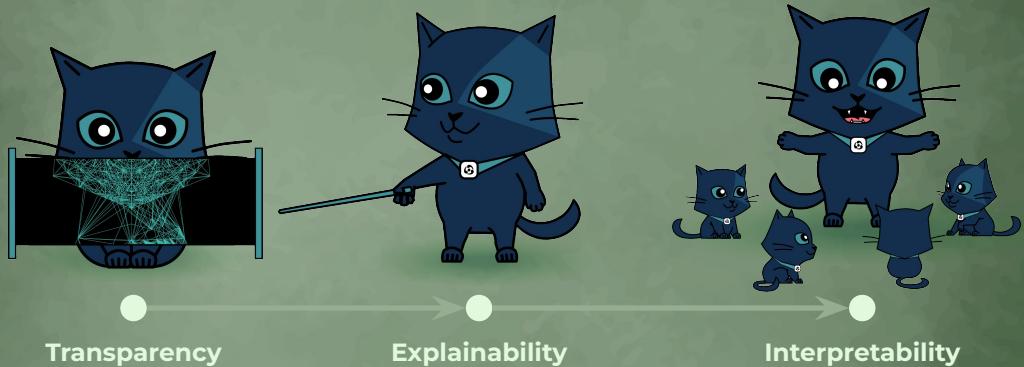
2.7 Transparency, Interpretability and Explainability

31

Let's cover the next set of requirements. They are **transparency, interpretability and explainability**.

Example. *Imagine an AI system that is designed to provide personalised product recommendations to online shoppers based on their browsing and purchase history. E.g., the AI system can exhibit **transparency, explainability** and **interpretability** characteristics by:*

- *Providing explanations for its product recommendations to the users. When suggesting a particular item, it accompanies the recommendation with clear and interpretable reasons. For instance, it might state that the recommendation is based on the user's previous purchase history, similar items purchased by other customers, or popular trends (and not because it's secretly gathering information about you while you are not paying attention)*
- *Allowing users to access and review their own data that the system uses to generate recommendations. These characteristics empower users to understand why certain recommendations are made and gain insights into their preferences.*
- *Providing clear information about the purpose of data collection, the scope of data usage, and how long the data will be retained.*



THE AIM OF XAI is for users to be able to understand how AI-based systems arrive at their results, thus **INCREASING USERS' TRUST** in them

Transparency: The level of visibility of the algorithm and data used by the AI-based system

Interpretability: The level of understanding how the underlying AI technology works

Sec. 2.7

Explainability: The level of understanding how the AI-based system came up with a given result

Explainable AI (XAI): The field of study related to understanding the factors that influence AI system outputs

32

Some of today's AI tools are able to produce highly-accurate results, but are also highly complex if not outright opaque, rendering their workings difficult to interpret. In the real world AI-based systems are typically applied in areas where users need to trust those systems. This could be for safety reasons, there could be demand for privacy or they might provide potentially life changing predictions and decisions. The complexity of AI-based systems has led to the field of "**Explainable AI**" (**XAI**). The aim of XAI is for users to be able to understand how AI-based systems arrive at their results, thus increasing users' trust in them.



REASONS FOR XAI according to THE ROYAL SOCIETY:

- giving users confidence in the system
- safeguarding against bias
- meeting regulatory standards or policy requirements
- improving system design
- assessing risk, robustness, and vulnerability
- understanding and verifying the outputs from a system
- autonomy, making the user feel empowered and meeting social value

Sec. 2.7

33

According to The Royal Society, there are several reasons for XAI, including:

- giving users confidence in the system;
- safeguarding against bias;
- meeting regulatory standards or policy requirements;
- improving system design;
- assessing risk, robustness, and vulnerability;
- understanding and verifying the outputs from a system;
- and autonomy, making the user feel empowered and meeting social value.

XAI has to be:

- **interpretable**, implying some sense of understanding how the technology works
- **explainable**, implying that users can understand why or how a conclusion was reached
- **transparent**, implying some level of accessibility to the data or algorithm
- **justifiable**, implying there is an understanding of the case in support of a particular outcome
- **contestable**, implying users have the information they need to argue against a decision

according to
ISTQB®

according to
THE ROYAL SOCIETY

Sec. 2.7

34

A range of terms describing some desired characteristics of an XAI can be found in this syllabus (the first three) and in the Royal Society's materials (all five). The explainable AI has to be:

- interpretable, implying some sense of understanding how the technology works;
- explainable, implying that users can understand why or how a conclusion was reached;
- transparent, implying some level of accessibility to the data or algorithm;
- justifiable, implying there is an understanding of the case in support of a particular outcome;
- contestable, implying users have the information they need to argue against a decision.

2.8 Safety and AI

And lastly we have **safety**,



Safety: The expectation that a system does not, under defined conditions, lead to a state in which human life, health, property, or the environment is endangered

Sec. 2.8

36

which implies that an AI-based system will not cause harm to people, property or the environment.

Example. Imagine an AI system that controls autonomous vehicles on public roads. **Safety** is a critical characteristic in this scenario, as the AI system must ensure the well-being of passengers, pedestrians, and other vehicles.

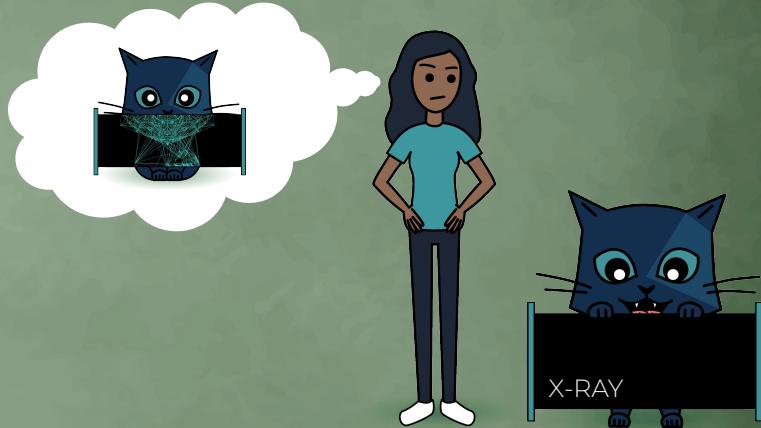
Collision Avoidance: The AI system is programmed to prioritise safety by actively monitoring the environment, detecting obstacles, and making decisions to avoid collisions. **Adhering to Traffic Rules:** The AI system follows traffic laws and regulations, including speed limits, traffic signals, and right-of-way rules. **Handling Emergency Situations:** The AI system is trained to handle emergency scenarios, such as sudden braking, swerving, or evasive manoeuvres, to avoid accidents. **Continual Monitoring and Redundancy:**

The AI system employs redundant systems and continual monitoring to enhance safety. It uses redundant sensors, backup control systems, and fail-safe mechanisms to minimise the impact of hardware or software failures.

Robust Testing and Validation: Extensive testing helps identify and address potential safety vulnerabilities before the system is deployed on public roads.

AI-BASED SYSTEM is often presented as a BLACK BOX with CHARACTERISTICS like:

- complexity
- non-determinism
- probabilistic nature
- self-learning
- lack of transparency
- lack of explainability
- lack of interpretability
- lack of robustness



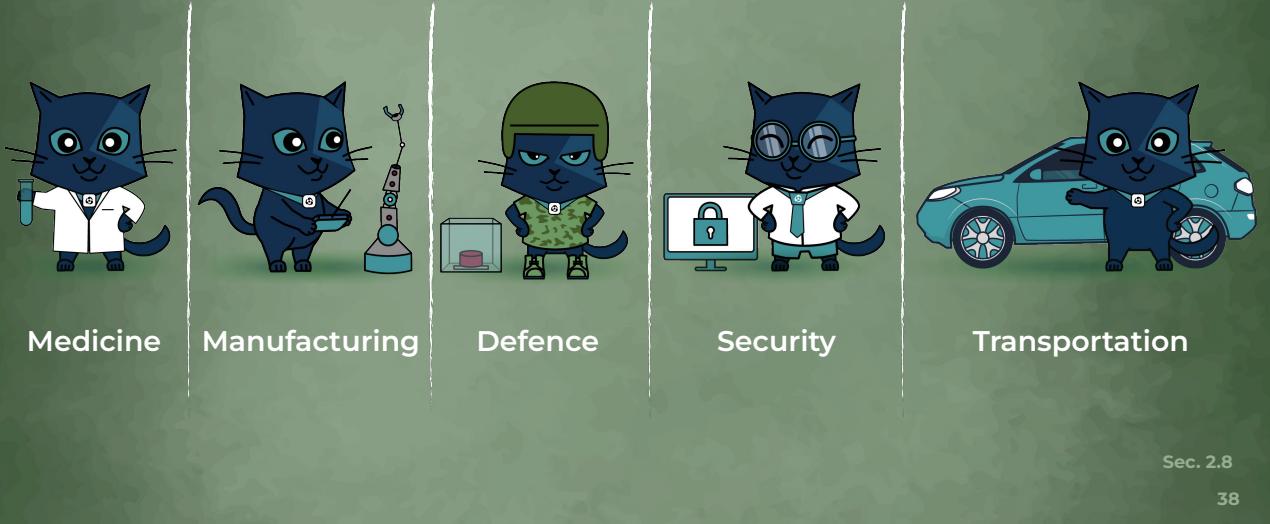
This makes it more difficult to trust an AI-based system and ensure that it is safe

Sec. 2.8

37

AI-based systems are often presented as black boxes with characteristics like complexity, non-determinism, probabilistic nature, self-learning, lack of transparency/interpretability/explainability as well as lack of robustness, which makes it more difficult to ensure they are safe.

AI-BASED SYSTEMS HAVE THE POTENTIAL TO AFFECT SAFETY



Sec. 2.8

38

For example, AI-based systems working in the fields of medicine, manufacturing, defence, security, and transportation have the potential to affect safety.

