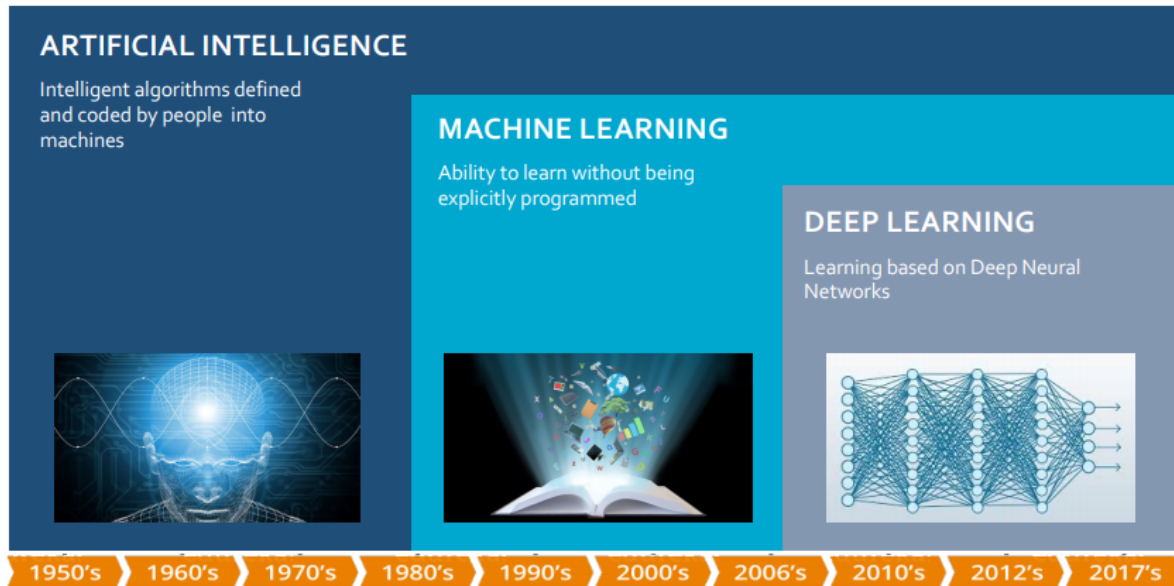


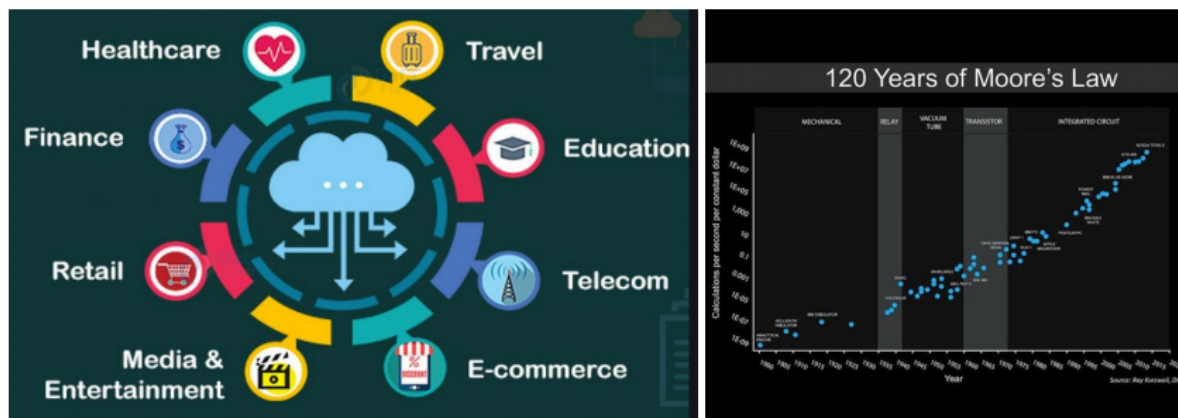
Lesson_14_IBM's_swag

A brief history of AI



Many of the AI ethics concerns are about the misuse of the technology, but some of them are related to specific techniques.

Data and computing power



The increased use of Machine learning techniques, which are the ones most related to AI ethics concerns, is due to the presence of a huge amount of data in different

areas and also the computing power to deal with that amount.

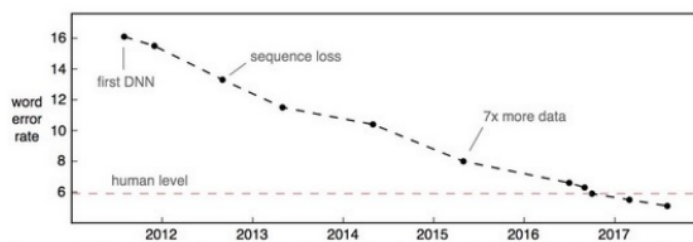
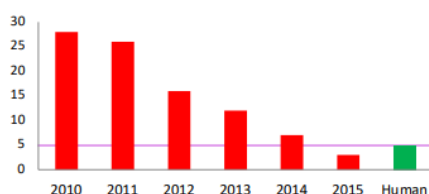
Image and natural language interpretation



Woman holding a cask of bananas



A group of young people playing frisbee



This led to a lot of success in many application where the AI system perceives the external world and then responds to it (e.g, image interpretation, speech recognition ...). Already few years ago machines were able to outperform humans in a specific task.

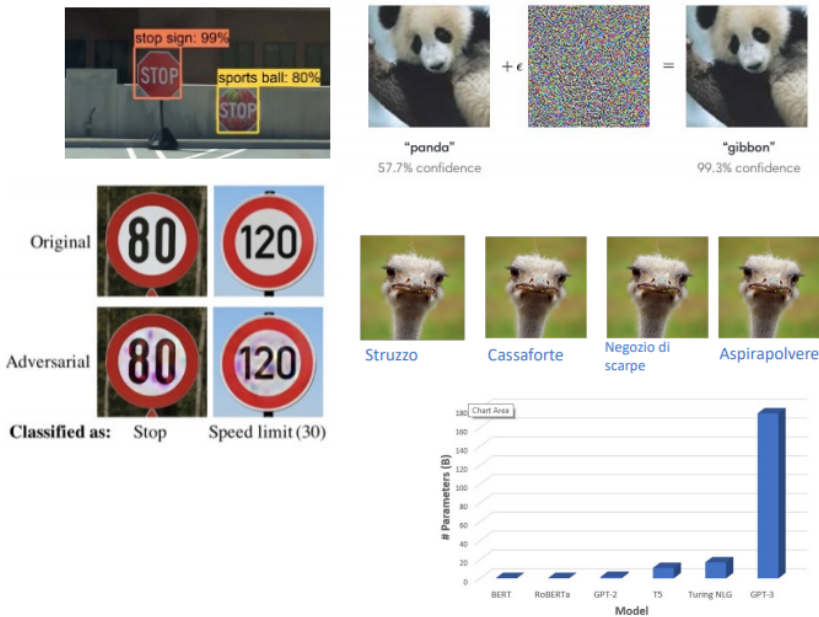
Some AI applications

West Point, April 2017

- Digital assistants:
 - Home assistants (Alexa)
 - Travel assistants (Waze)
- Driving/travel support:
 - Auto-pilot (Tesla)
 - Ride-sharing apps (Uber, Lyft)
- Customer care:
 - Client service chatbots
- Online recommendations:
 - Friend recommendations (Facebook)
 - Purchase recommendations (Amazon)
 - Movie recommendations (Netflix)
- Media and news:
 - Ad placement (Google)
 - News curation
- Healthcare:
 - Medical image analysis
 - Treatment plan recommendation
- Financial services:
 - Credit risk scoring
 - Loan approval
 - Fraud detection
- Job market:
 - Resume prioritization
- Judicial system:
 - Recidivism prediction (Compas)



Many applications derive from the machines' capabilities.



AI limitations

- Narrow AI
 - Solves well specific problems
- Lack of robustness and adaptability
- Needs a lot of resources
 - Data and computing power



AI has a lot of capabilities, but also a lot of limitations. The lack of robustness and adaptability are related: AI focused on a specific problem lacks the capability of transferring the learned knowledge to another problem.

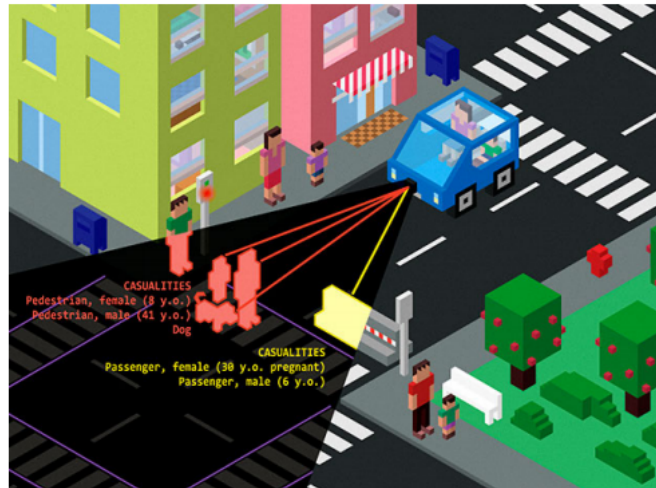
The errors in image classification can be due to noise or other external interventions, which would not cause trouble to humans, are another sign of the lack of robustness. These are the reasons why there is reluctance of using these models where some danger may derive from these errors or there's need of accountability.

Ethical issues -- examples



Ethical issues related to AI in some companies

Can we trust AI's decisions?



Usual Trolley Problem with self-driving cars: harm pedestrians or the passengers?

On what values should be the car programmed on for us to be able to trust the decisions it will make?

AI Ethics



Multidisciplinary field of study



How to optimize AI's beneficial impact while reducing risks and adverse outcomes



How to design and build AI systems that are aware of the values and principles to be followed in the deployment scenarios



To identify, study, and propose technical and nontechnical solutions for ethics issues arising from the pervasive use of AI in life and society

AI ethics is a multidisciplinary field of study that involves not only AI experts but other experts who are able to assess the **impact** that the new technologies can have on society (sociologist, philosophers, economist...). Also policy makers, who choose the regulations of our every-day life.

The goals of this field of study are:

- **Optimize** AI's beneficial impact, while reducing risks.

- Build machines that will work according to some values **relevant** for their deployment scenario.
- Development of guidelines, best-practices, standards, regulations and laws; these can help improving the technical solutions and, maybe, completing them.



AI needs raise ethical problems:

- **Data privacy & Governance:** the gathering of data involves personal and sensitive data, what happens to them? Who gets them?
- **Explainability & Transparency:** how did the machine came up with a certain conclusion is hardly explainable. Also, they're not very transparent: how they been built? What design and principles have been followed?
They are kind of a black box and this generates concerns and mistrust.
- **Fairness & alignment:** the AI systems can make recomendations, but how can we be sure that they are aware of human values or aligned to them? We want the decisions to be **fair** and not create discrimination in the same group of people, for example.
- **Accountability:** we know that the machine learning systems are not perfect and will commit some errors, like humans do. But when a machine makes a mistake who's to be held accountable?
- **Profiling & Manipulation:** profiling is done to infer an individual's preferences, in order to show advertisment that we are interested in. If our preferences are variegated, it's hard to understand what we may like or not. But if they are polarized, then is easier. Providing us content that may steer our preferences towards polarization is a form of manipulation.

- **AI pervasiveness:** the high pervasiveness of AI implies that a misuse of the technology has a very large impact across different application domains, and also brings many transformations (jobs, society ...). Are we keeping track of these transformations, are the possibly negative aspects interpreted so that we can put in place some solutions?

AI is not a neutral technology

- Misuse must be avoided
- But AI needs to be designed and developed with the right properties
 - Fair, explainable, robust, ...

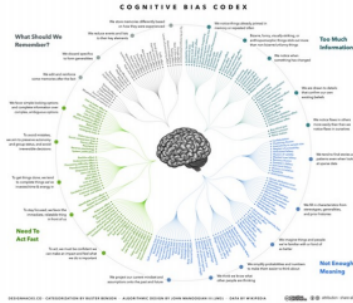


Misuse must be avoided of course, but AI technologies must be developed with the right properties in order to be sure that the way we deal with AI is ethical and responsible.

AI is not a neutral technology: the way it is designed and developed and the properties that are embedded into it are also part of how the technology will behave, It's not just a matter of avoiding misuse.

AI fairness

- Bias: prejudice for or against something
- As a consequence of bias, one could behave unfairly to certain groups compared to others
- Why should AI be biased?
 - Trained on data provided by people, and people are biased



West Point, April 2021



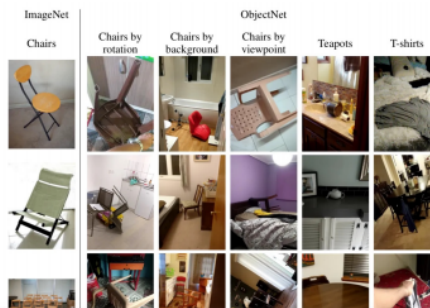
Very similar to the logic in the Fairness in algorithmic Decision-making

Why should be AI biased? Biases lead to discriminatory behaviour

People who build the AI will also generate the data on which it trains. But people are biased... Any step involving humans during the development will introduce some biases.

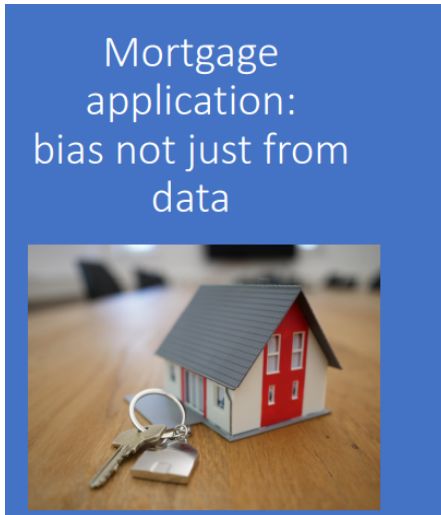
AI bias: ImageNet

- 14M images, used to train image interpretation AI systems
- Bias in the data distribution and in the data labels (Mturk people)



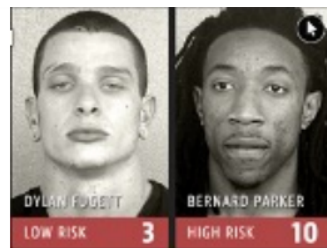
West Point, April 2021





- Training data
 - Ex. : correlation gender-acceptance
- Design decisions:
 - Ex.: prioritized motivations for loan applications
 - Buying a house
 - Paying school fees
 - Paying legal fees
 - Loan applications with these motivations are prioritized
 - If one of them is omitted, the relevant community will be penalized

Biases come from not only from the data but also from every decision we make. For example, if we build a mortgage application acceptance AI system: given a mortgage application it decides whether it should be accepted or declined. During the development of the system, we might want to introduce some prioritized motivations for loan applications. Given these priorities, we assume to know what is best for most people. But what might be best for one ethnicity might not be for another one. Therefore, these priorities introduce a form of bias.



- Overall accuracy is the same, regardless of race (**overall accuracy equality**)
- Likelihood of recidivism among defendants labeled as medium or high risk is similar, regardless of race (**predictive parity**)
- But ... false positive and false negative rates are very different

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

West Point, April 2021



COMPAS System again

The algorithm's output is a score that describes the probability that the subject will commit a crime again. The judges would use these scores to decide if the subject should be staying in jail while waiting for the final judgement.

The algorithm has been found to be biased, but it highlighted the fact that there are more definition of fairness: an algorithm can be considered fair according to one notion but not to another one. In this case, the accuracy of the algorithm over the two groups (black & white) was more or less the same: there was equality in overall accuracy. According to this definition of fairness, the algorithm was ok. But looking at the errors, the false positives (high risk when it should be low) and false negative (low when high), it was found that the rates were very different, when they should have been similar at least. According to this definition of fairness, the algorithm was not fair.

This shows that to be able to assess an AI system fairness, we must choose the notion **relevant** for that application domain.

Many decision points

- **Individual vs group fairness:**
 - similar individuals should receive similar treatments or outcomes, vs
 - groups defined by protected attributes should receive similar treatments or outcomes
- **Context-dependent definition(s) of fairness**
- **Acceptable bias threshold**
- **When to detect bias:**
 - training data or learned model

Source: *Fairness and Machine Learning* by Solon Barocas, Moritz Hardt, Arvind Narayanan (<https://www.fairmlbook.org>)

Name	Closest relative	Note	Reference
Statistical parity	Independence	Equivalent	Dwork et al. (2011)
Group fairness	Independence	Equivalent	
Demographic parity	Independence	Equivalent	
Conditional statistical parity	Independence	Relaxation	Corbett-Davies et al. (2017)
Darlington criterion (4)	Independence	Equivalent	Darlington (1971)
Equal opportunity	Separation	Relaxation	Hardt, Price, Srebro (2016)
Equalized odds	Separation	Equivalent	Hardt, Price, Srebro (2016)
Conditional procedure accuracy	Separation	Equivalent	Berk et al. (2017)
Avoiding disparate mistreatment	Separation	Equivalent	Zafar et al. (2017)
Balance for the negative class	Separation	Relaxation	Kleinberg, Mullainathan, Raghavan (2016)
Balance for the positive class	Separation	Relaxation	Kleinberg, Mullainathan, Raghavan (2016)
Predictive equality	Separation	Relaxation	Chouldechova (2016)
Equalized correlations	Separation	Relaxation	Woodworth (2017)
Darlington criterion (3)	Separation	Relaxation	Darlington (1971)
Cleary model	Sufficiency	Equivalent	Cleary (1966)
Conditional use accuracy	Sufficiency	Equivalent	Berk et al. (2017)
Predictive parity	Sufficiency	Relaxation	Chouldechova (2016)
Calibration within groups	Sufficiency	Equivalent	Chouldechova (2016)
Darlington criterion (1), (2)	Sufficiency	Relaxation	Darlington (1971)

West Point, April 2021



Once we have chosen the definition of fairness, we must find an acceptable threshold for bias since it's unavoidable (80% is the typical one).

The developers should be educated to recognize their own biases and understand what biases they can inject into the system. Another method to keep biases in check is to form heterogenous developing teams, able to spot each others biases.

AI
explainability:
AI systems
cannot be
black boxes

The **General Data Protection Regulation (GDPR)**

- Limits to **decision-making** based solely on **automated processing** and profiling (Art.22)
- Right to be provided with **meaningful information** about the **logic** involved in the decision (Art.13 (2) f. and 15 (1) h)

The data subjects that give the data to the AI system have the **right** to have an explanation of the decisions made by the system.

Again, a system that does not give any explanation for its outcomes it's unlikely to be trusted, in general.

A system that works in a certain domain should be able to generate different explanation for different audiences. For example, a system working in the medical domain should be able to generate explanation for the doctors, the patients and so on, in a way that is understandable for them.

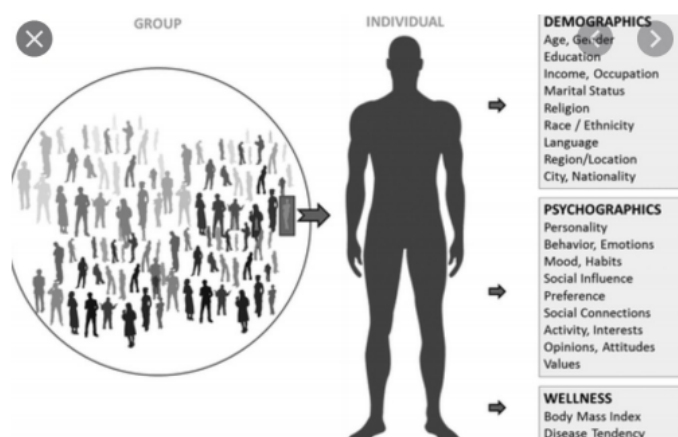
Data handling: the General Data Protection Regulation (GDPR)



The GDPR is about data protection and sets a standard and gives regulation, very important for EU. Many countries lack a comprehensive document like this one, even if it is partially used in some other regions of the world.

Profiling and manipulation

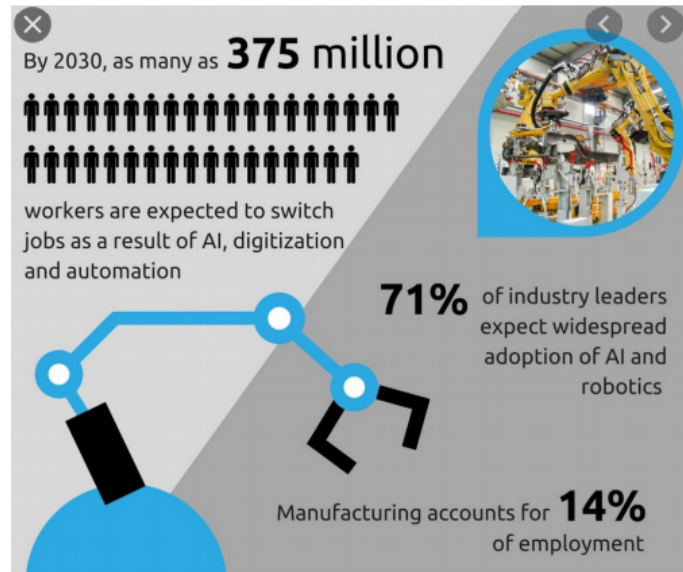
- From actions to profiles
 - Like, text, images, follow, ...
- AI can infer our preferences, and use them to advertise products that we probably like
 - Easier if our preferences are bipolar



As we said before, our preferences are tried to be made more polarized by the advertising systems because it eases the job of giving highly preferred ads.

Impact on the workforce

- Many jobs will disappear, and many others will be created
- All jobs will change

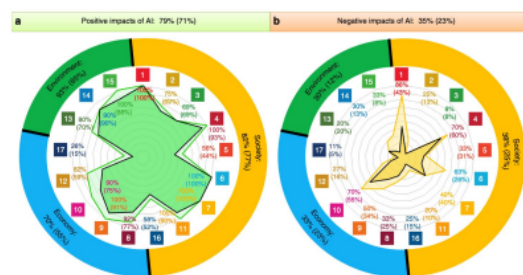


People have to be convinced of the transition benefits because it's easier to imagine existing jobs disappear rather than imagine those which will be created.

It is expected that AI will change the jobs, for now it will be a support for the human professional figures but also will guide robotics in the manufacturing environment.

A vision of the future (2030)

- 17 goals, 169 targets
- Very difficult path
 - The pandemic has worsened the situation
- AI can help in achieving the SDGs
- COVID: vaccines in less than one year!



The UN have put together a comprehensive vision of what it means to improve the world. They used 17 goals that define how the world should look like in 2030. A recent study showed how much the current AI techniques have been used to move towards these SDG and how much it has brought further away from them.

This SDG can help to steer the use of technology in a desirable direction.

Where was IBM when the AI fell?

IBM, technology, and AI

- 110 years
- Hardware e software
- Enterprise AI: AI solutions for other companies
 - Banks and financial institutions
 - Governments
 - Aeroports
 - Hospitals
 - ...



Summit, IBM



Chess: IBM Deep Blue, 1997



Quantum computer, IBM



Jeopardy: IBM Watson, 2011



Project Debater, 2020

Overview of the IBM technology in the AI field, interesting but doesn't seem to be strictly related to ethics or whatever else relevant for the course.

IBM Principles of Trust and Transparency (2017)



The purpose of AI is to **augment** human intelligence



Data and insights belong to their creator



New technology, including AI systems, must be **transparent** and **explainable**

IBM in terms of ethics defined in 2017 some high-level principles:

- AI should *augment* human intelligence rather than replace it, therefore helping humans in their decisions.

- Data gathered from clients belong to them: they will not be reused to improve solutions for other clients.
- The technology should be transparent, for example a human being should always know whether s/he is interacting with a human or an AI system. The technology should be explainable



AI PRINCIPLES in the world – a comprehensive view

Actors:

- Private sector
- Inter-governmental
- Multistakeholder
- Governments
- Civil society

Main themes:

- Human rights
- Human values
- Responsibility
- Human control
- Fairness
- Transparency and explainability
- Safety and Security
- Accountability
- Privacy

Principled AI Project, Berkman Klein's Cyberlaw Clinic, 2019



Each ray represents one set of principles while the colors represent actors

These principles were the starting point of a long path and from the principles now we should take some concrete actions, implementing them in the technology.

What does it mean to TRUST a decision made by a machine? (Other than it is accurate and respect privacy)

West Point, April 2021



Is it **fair**, or is it going to make discriminatory decisions?



Is it possible to understand **why** it made that decision, or is it a black box?



Is it **robust**?



Is it **transparent**?



Trust definition using these four pillars listed above

Everyday Ethics for Artificial Intelligence

AI Fairness 360

This extensible open source toolkit can help you examine, report, and mitigate discrimination and bias in machine learning models throughout the AI application lifecycle. We invite you to use and improve it.

Python API Docs | Get Python Code | Get R Code

West Point, April 2021

- **Technical solutions to detect and mitigate AI bias**
 - Research work
 - Watson OpenScale
 - Open-source libraries: AI fairness 360
- **Developers' education and training**
 - AI bias education modules for all IBMers
 - Developers' awareness material
 - Revised methodologies for the AI pipeline
 - Adoption strategies
 - Governance frameworks
 - Consultations with all stakeholders
 - Design thinking sessions



IBM's initiatives to respect the defined principles

Other than tools to inspect and detect unfair behaviours, IBM investend into educating its developers to help them undstand those tools, when to used them and why.

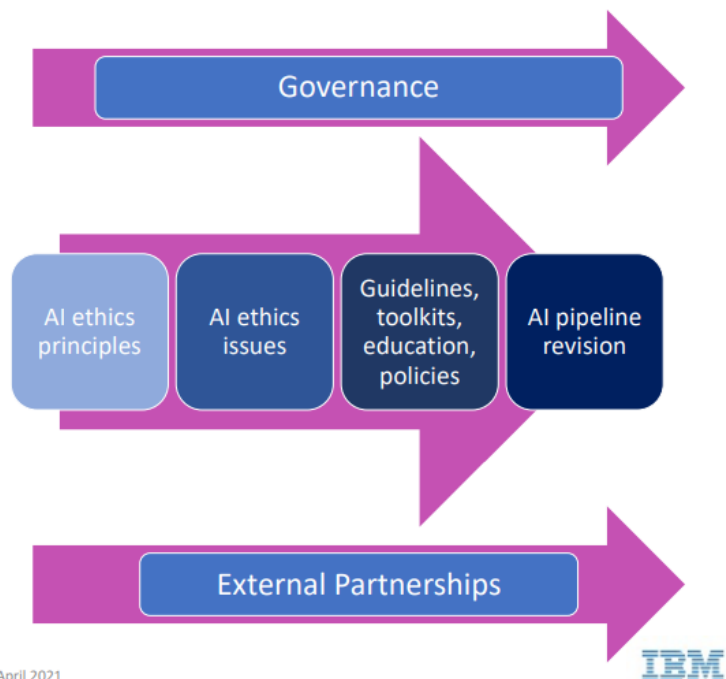


- AI factsheet
 - Transparency by documentation
 - Design a development choices
 - Not just a checklist
 - Self-assessment and beyond
- Useful to
 - Developers
 - Clients
 - Users regulators/auditors
- Aligned with EC High Level Expert Group on AI self-assessment list (ALTAI)
- AI factsheet 360

IBM wants to be transparent on all the design decisions that it has made during the development process and to do that uses AI factsheet, which are an implementation that follows the idea of **transparency by documentation**.

AI factsheet 360 open-source toolkit to understand this concept and possibly use them.

From principles to practice: a multi-dimensional space



West Point, April 2021

It is a journey, it takes a long time to change how things are done, especially in big companies.

Governance: the IBM AI Ethics board

- Mission
 - Awareness and coordination
 - Internal education and retraining
 - Linking research to services and platforms
 - Advice to business units
 - Internal governance framework
 - Define policies and advice regulators
- Risk-based approach for the BUs
 - Vetting based on three dimensions (tech, use, client)



This board supervises all the mentioned activities and it decides whether if an AI solution can be offered to a certain client, if there are ethical uncertainties about it. The board takes that decision according to certain parameters (e.g, Bias threshold), considering what the deployment the solution will have and its context and the client itself.

Partnerships

Academia
Companies
Governments
Civil society
organizations

Multi-disciplinary and
multi-stakeholder

Asilomar AI principles

RESEARCH	ETHICS AND VALUES	EMERGING ISSUES
1. Research goals	6. Safety	19. Quality of training
2. Research funding	7. Public transparency	20. Interpretability
3. Science-policy links	8. Auditing transparency	21. Bias
4. Research culture	9. Responsibility	22. Resilience
5. Share resources	10. Value alignment	23. Robustness
	11. Human values	24. Self-improvement
	12. Personal privacy	25. Control/guard
	13. Liberty and privacy	
	14. Global benefit	
	15. Global prosperity	
	16. Human control	
	17. Non-maleficence	
	18. AI arms race	

IBM | MIT

AAAI / ACM conference on ARTIFICIAL INTELLIGENCE, ETHICS, AND SOCIETY

AI for Good Global Summit
An ITU experience

Version II - For Public Discussion | **IEEE**
Advancing Technology for Humanity

ETHICALLY ALIGNED DESIGN
A Guide for Promoting Human Well-being with Transdisciplinary and Interdisciplinary Approaches

Partnership on AI
to benefit people and society

One organization

Innovate and share the best practices for using and developing AI technologies and providing a global platform to discuss how AI will influence people and society.

7 Thematic Pillars

AI Labor and the Economy
AI Transparency and Accountability
Governance, Regulation, Policy and AI Systems
Human and Machine Influence of AI
Human and Machine Influence of AI

Logos: Microsoft, Google, Facebook, Amazon, Intel, SAP, Oracle, IBM, etc.

WORLD ECONOMIC FORUM

West Point, April 2021

EUROPEAN COMMISSION
Ethics Guidelines for Trustworthy AI

EUROPEAN COMMISSION
Policy and Regulatory Recommendations for Trustworthy AI

Partnerships with external entity is necessary to better understand the issues: it has to be multi-disciplinary, but also multi-stakeholder because all the different voices have to be there to understand the impact on society.

Not just AI

- Neurotechnologies
 - Huge potential for healthcare
 - Reading/writing neurodata
 - Additional issues around privacy, agency, and identity
- Quantum computing
 - How to responsibly use such a huge computing power?



There are other fields that can be combined with AI and raise ethical concerns: one example is **neurotechnologies**; what would happen to the concept of privacy if we were able to read neurodata (nervous system's signals)?

Quantum computing will provide huge computational power and some protection based on complexity may fall: how can we protect it from cyber attacks? How can we ensure that this technology will be used in a responsible way?

List of useful links in the slides.