# AI Accountability

SOULEIMA ZGHAB

# The Rise of Algorithms and the Need for Countervailing Checks

- The rise of Artificial Intelligence

- Significant failings across a range of applications

# The Rise of Algorithms and the Need for Cou...

- The rise
- Significa

## Amazon ditched AI recruiting tool that favored men for technical jobs

**Specialists had been building computer programs since 2014 to review résumés in an effort to automate the search process**

# The Ri... or Counte...

- The rise of Art...
- Significant faili...

# The Rise of Algorithms and the Need for Countervailing Checks

- The rise of Artificial Intelligence

- Significant failings across a range of applications

It is increasingly important that we identify the best ways to keep them **accountable**.

# Key concepts

Accountability - "The state of being held responsible or answerable for a system, its behavior and its potential impacts" [Raji et al]

# Who should be held accountable ?

- Although algorithms themselves cannot be held accountable as they are not moral or legal agents ?

# Who should be held accountable ?

- Although algorithms themselves cannot be held accountable as they are not moral or legal agents ?


THE MIC IS YOURS
memegenerator.net

# Who should be held accountable ?

- Although algorithms themselves cannot be held accountable as they are not moral or legal agents

> "Nearly 50% of the surveyed developers believe that the humans creating AI should be responsible for considering the ramifications of the technology. Not the bosses. Not the middle managers. The coders."
>
> – Mark Wilson, Fast Company on Stack Overflow's Developer Survey Results 2018

# Who should be held accountable ?

- Although algorithms themselves cannot be held accountable as they are not moral or legal agents,

- The organizations designing and deploying algorithms should be held accountable  through governance structures.

> "Nearly 50% of the surveyed developers believe that the humans creating AI should be responsible for considering the ramifications of the technology. Not the bosses. Not the middle managers. The coders."
>
> – Mark Wilson, Fast Company on Stack Overflow's Developer Survey Results 2018

# Key concepts

Accountability - "The state of being held responsible or answerable for a system, its behavior and its potential impacts" [Raji et al]

Governance - "The system by which the whole organization is directed, controlled, and held accountable to achieve its core purpose over the long term" [ISO 37000 standard]

# Key concepts

Accountability - "The state of being held responsible or answerable for a system, its behavior and its potential impacts" [Raji et al]

Governance - "The system by which the whole organization is directed, controlled, and held accountable to achieve its core purpose over the long term" [ISO 37000 standard]

Auditing - "An Independent evaluation of conformance of software products and processes to applicable regulations, standards, guidelines, plans, specifications, and procedures." [IEEE standard]

# Accountability: A Conceptual Frame

- Accountability is the foundation of <u>trust</u> in society.

- Accountability implies <u>an obligation to report and justify algorithmic decision-making,</u> and to mitigate any negative social impacts or potential harms.

- Algorithmic accountability vs Algorithmic justice ?

# Accountability: A Conceptual Frame

- Accountability is the foundation of trust in society.

- Accountability implies <u>an obligation to report and justify algorithmic decision-making,</u> and to mitigate any negative social impacts or potential harms.

- Algorithmic accountability vs Algorithmic justice ?

# Accountability: A Conceptual Frame

- Accountability is the foundation of trust in society.

- Accountability implies <u>an obligation to report and justify algorithmic decision-making,</u> and to mitigate any negative social impacts or potential harms.

- Algorithmic accountability vs Algorithmic justice ?
  - Algorithmic accountability: the responsibility of algorithm designers to provide evidence of potential or realized harms,
  - Algorithmic justice: the ability to provide redress from harms.

# Accountability related to AI

- Three "senses" of accountability related to AI exist in the literature :
  - Accountability is a feature of the AI system itself
  - Who is most responsible for what effect within the sociotechnical system
  - A feature of the broader sociotechnical system that develops, procures, deploys and uses AI
    - AI Now proposes an Algorithmic Impact Assessment framework as a means of building accountability into the broader sociotechnical system in which AI is deployed, only part of which would include responsibility determinations.

# Acc

Government Gouvernement
of Canada du Canada

**Algorithmic Impact Assessment**

Home > Open Government

## Algorithmic Impact Assessment

ℹ Information in the AIA is only stored locally on your computer, and the Government of Canada does not have
access to the information you place into the tool. If you wish to keep your work, please save the data locally for
future use by using the 'Save' button. You can reload a previously saved AIA form using the 'Upload JSON File' button.

Upload JSON File

### Algorithmic Impact Assessment v0.9

#### 1. What is the Algorithmic Impact Assessment?

The AIA is a questionnaire designed to help you assess and mitigate the impacts associated with deploying an automated decision system. The
AIA also helps identify the impact level of your automated decision system under the Directive on Automated Decision-Making. The questions

✔ Impact Level: 1          Current Score: 0          Raw Impact Score: 0          Mitigation Score: 0

- Three "s
  - Accoun
  - Who is
  - A featu
    - AI Now
      sociote

# General Data Protection Regulation (GDPR)

The General Data Protection Regulation (GDPR) integrates accountability as a principle which requires that organisations put in place appropriate technical and organisational measures and be able to demonstrate what they did and its effectiveness when requested.

1. Personal data shall be:

   (a) processed lawfully, fairly and in a transparent manner in relation to the data subject ('lawfulness, fairness and transparency');

   (b) collected for specified, explicit and legitimate purposes and not further processed in a manner that is incompatible with those purposes; further processing for archiving purposes in the public interest, scientific or historical research purposes or statistical purposes shall, in accordance with Article 89(1), not be considered to be incompatible with the initial purposes ('purpose limitation');

   (c) adequate, relevant and limited to what is necessary in relation to the purposes for which they are processed ('data minimisation');

   (d) accurate and, where necessary, kept up to date; every reasonable step must be taken to

   protection against unauthorised or unlawful processing and against accidental loss, destruction or damage, using appropriate technical or organisational measures ('integrity and confidentiality').

2. The controller shall be responsible for, and be able to demonstrate compliance with, paragraph 1 ('accountability').

# General Data Protection Regulation (GDPR)

The General Data Protection Regulation (GDPR) integrates accountability as a principle which requires that organisations put in place appropriate technical and organisational measures and be able to demonstrate what they did and its effectiveness when requested.

Under the General Data Protection Regulation (GDPR), the principle of accountability is intrinsically linked to the principle of **transparency**.

Transparency empowers data subjects to hold data controllers and processors accountable and to exercise control over their personal data. Accountability requires transparency of processing operations, however transparency does not constitute accountability [3].

# Institutional Mechanisms

- "Institutional mechanisms" are processes that shape or clarify the incentives of the people involved in AI development, make their behavior more transparent, or enable accountability for their behavior.

# Institutional Mechanisms

- "Institutional mechanisms" are processes that shape or clarify the incentives of the people involved in AI development, make their behavior more transparent, or enable accountability for their behavior.

- Institutional mechanisms help to ensure that individuals or organizations making claims regarding AI development are incentivized to be diligent in developing AI responsibly and that other stakeholders can verify that behavior.

# Institutional Mechanisms

- Institutions can **shape** incentives or constrain behavior in various ways.
- Institutional mechanisms **can help clarify an organization's goals and values**, which in turn can provide a basis for evaluating their claims.
- Institutional mechanisms **can increase transparency** regarding an organization's AI development processes in order to permit others to more easily verify compliance with appropriate norms, regulations, or agreements.
- Institutional mechanisms can **create incentives** for organizations to act in ways that are responsible
- Institutional mechanisms can **foster exchange of information** between developers

Mechanism with the potential for improving the **verifiability of claims** in AI development: third party auditing, red team exercises, bias and safety bounties, and sharing of AI incidents.

# Institutional Recommendations

- <u>The mechanism:</u> <span style="color:green">Third Party Auditing</span>

- <u>The problem:</u> The process of AI development is often opaque to those outside a given organization, and various barriers make it challenging for third parties to verify the claims being made by a developer. As a result, claims about system attributes may not be easily verified.

- <u>The recommendation:</u> A coalition of stakeholders should create a task force to research options for conducting and funding third party auditing of AI systems.

# Institutional Recommendations

- **The mechanism**: Red Team Exercises

- **The problem:** It is difficult for AI developers to address the "unknown unknowns" associated with AI systems, including limitations and risks that might be exploited by malicious actors. Further, existing red teaming approaches are insufficient for addressing these concerns in the AI context.

- **The recommendation:** Organizations developing AI should run red teaming exercises to explore risks associated with systems they develop, and should share best practices and tools for doing so.

# Institutional Recommendations

- **The mechanism:** Bias and Safety Bounties

- **The problem:** There is too little incentive, and no formal process, for individuals unaffiliated with a particular AI developer to seek out and report problems of AI bias and safety. As a result, broad-based scrutiny of AI systems for these properties is relatively rare.

- **The recommendation:** AI developers should pilot bias and safety bounties for AI systems to strengthen incentives and processes for broad-based scrutiny of AI systems.

# Institutional Recommendations

- **The mechanism:** Sharing of AI Incidents

- **The problem:** Claims about AI systems can be scrutinized more effectively if there is common knowledge of the potential risks of such systems. However, cases of desired or unexpected behavior by AI systems are infrequently shared since it is costly to do unilaterally.

- **The recommendation:** AI developers should share more information about AI incidents, including through collaborative channels.

# Software Mechanisms

Software mechanisms involve **shaping** and **revealing** the functionality of existing AI systems.

They can support verification of new types of claims or verify existing claims with higher confidence.

For example, an AI developer that wants to provide evidence for the claim that "user data is kept private" can help build trust in the lab's compliance with a formal framework such as differential privacy, but non-experts may have in mind a different definition of privacy.

Enabling verification of claims about system:

- Reproducibility of technical results in AI

- Formal verification

- The empirical verification and validation of machine learning by machine learning

- Practical verification is the use of scientific protocols to characterize a model's data, assumptions, and performance

# Software Recommendations

- **The mechanism:** Audit Trails

- **The problem:** AI systems **lack traceable logs** of steps taken in problem-definition, design, development, and operation, leading to a lack of accountability for subsequent claims about those systems' properties and impacts.

- **The recommendation:** Standards setting bodies should work with academia and industry to develop audit trail requirements for safety-critical applications of AI systems.

# Software Recommendations

- <u>The mechanism:</u> Interpretability

- <u>The problem:</u> It's difficult to verify claims about **"black-box"** AI systems that make predictions without explanations or visibility into their inner workings. This problem is compounded by a lack of consensus on what interpretability means.

- <u>The recommendation:</u> Organizations developing AI and funding bodies should support research into the interpretability of AI systems, with a focus on supporting risk assessment and auditing.

# Software Recommendations

- **The mechanism:** Privacy-Preserving Machine Learning

- **The problem:** A range of methods can potentially be used to verifiably safeguard the data and models involved in AI development. However, standards are lacking for evaluating new privacy-preserving machine learning techniques, and the ability to implement them currently lies outside a typical AI developer's skill set.

- **The recommendation:** AI developers should develop, share, and use suites of tools for privacy preserving machine learning that include measures of performance against common standards.

# Discussion

- Which mechanism do you think is the most urgent ?

- Is there another aspect to be taken into consideration ?

# References

[1] https://www.ic.gc.ca/eic/site/133.nsf/vwapj/3_Discussion_Paper_-_Accountability_in_AI_EN.pdf/$FILE/3_Discussion_Paper_-_Accountability_in_AI_EN.pdf1

[2] https://datasciencemilan.medium.com/general-data-protection-regulation-gdpr-a-data-science-perspective-ecc02f6d6874

[3] Raji, Inioluwa Deborah, et al. "Closing the AI accountability gap: Defining an end-to-end framework for internal algorithmic auditing." Proceedings of the 2020 conference on fairness, accountability, and transparency. 2020.

[4] Brundage, Miles, et al. "Toward trustworthy AI development: mechanisms for supporting verifiable claims." arXiv preprint arXiv:2004.07213 (2020).

[5] https://www.theguardian.com/technology/2018/oct/10/amazon-hiring-ai-gender-bias-recruiting-engine

[6] https://www.dataprotection.ie/en/organisations/know-your-obligations/accountability-obligation