

Response to the European Commission's Proposed AI Act

6 AUGUST 2021

FACEBOOK

Table of Contents

Table of Contents	1
I. Introduction.....	1
II. Facebook is Committed to Responsible AI.....	2
III. The AI Act Must Allow for Refinement in Areas Where Foundational Questions Remain	3
A. Many of the AI Act’s Requirements Are Not Feasible or Assume Best Practices or Standards That Do Not Yet Exist	4
1. <i>The scope of many provisions reaches technologies that are not AI or not high-risk</i>	<i>4</i>
2. <i>The proposal’s attempt to promote fairness does not reflect the totality of causes of—and mitigations for—bias.....</i>	<i>8</i>
IV. We Need Collaboration to Develop Consensus Standards and Best Practices	11
A. Imposing Obligations Without Clear Standards Risks Serious Costs to Innovation.....	11
1. <i>What constitutes reasonable efforts to detect and mitigate bias in AI systems?</i>	<i>13</i>
2. <i>What are the appropriate standards and measures for robustness, accuracy, security, risk management and assessment?</i>	<i>14</i>
3. <i>How do we reconcile privacy with the need to measure bias using special categories of personal data?</i>	<i>15</i>
4. <i>How can we encourage the creation of more high-quality, open data sets?</i>	<i>16</i>
5. <i>How can we encourage greater data sharing across industry and between the public and private sectors?</i>	<i>16</i>
B. Advancing the EU’s Digital Single Market Strategy Requires Union-wide Consistency in Regulation and Enforcement.....	17
V. Conclusion	17
VI. Endnotes	19

I. Introduction

Facebook Ireland (“Facebook”) welcomes the opportunity to provide comments on the Commission’s proposed Artificial Intelligence Act (the “AI Act”). AI is a uniquely powerful technology that must be used responsibly, and we believe thoughtful regulation can help ensure this responsible use.¹ For that reason, Facebook looks forward to working with policymakers in Europe to strengthen the Commission's proposal.

Facebook is a global leader in the creation and application of AI technology, and we understand first-hand that these are still early days for both AI technology and AI governance. The AI technologies that will transform Europe and the world are still being developed. In parallel, our understanding of the best approaches to govern those technologies are also still being developed. Over the past few years, governments, companies, and civil society organisations have all developed high-level principles for responsible AI development, deployment, and use.² For example, the European Commission’s High-Level Expert Group on AI developed the Ethics Guidelines for Trustworthy AI,³ and the OECD’s AI Governance Expert Group, including active involvement from Facebook, created the Principles on Artificial Intelligence,⁴ which have since been adopted by the G20.

These frameworks have provided a good foundation, but they also have left key questions unanswered. For example, although nearly every high-level framework discusses the importance of ensuring that AI systems are fair and operate in non-discriminatory ways, there are no standards or best practices to guide how to balance competing values such as data minimisation versus fairness (which may necessitate collecting data to measure and mitigate potential bias), or how to balance the desire for accurate data with a desire to avoid perpetuating existing biases that may be reflected in that data.

Imperfect knowledge, however, cannot be an excuse for inaction, and the Commission is right to begin this conversation now. The Commission can ensure that the greatest risks from specific uses of AI are contained, while allowing innovation to flourish in Europe. But to do so, the regulations must reflect the realities of how much we still do not know. As explained in more detail below, the Commission's proposal—while promising in many respects—often presumes that these questions do, in fact, have answers, when this is often not the case. We offer our comments on how the regulation can better reflect the current state of play around AI policy issues, while building in room for future refinement of these issues. We believe that society will benefit when AI technologies are trusted, safe, fair, accountable, and transparent. To achieve this, we need AI regulations that are clear, practicable, effective, and flexible. We look forward to working with European policymakers to ensure that the resulting regulation is clear in its expectations,

practicable, and technically feasible, so that it can unlock a new era of responsible AI innovation in Europe.

II. Facebook is Committed to Responsible AI

At Facebook, AI is an essential part of how we show people the content that matters to them and how we keep them safe when using our services. Facebook's AI models perform trillions of operations every day for the billions of people that use our technologies.⁵ For each of the more than 2 billion people who use our services each day, there could be more than a thousand posts that could appear in that person's feed at any given moment. AI helps us show the posts that are likely to be most meaningful to each user.⁶ Additionally, we use AI to help users discover new communities and content they may be interested in through tools like Pages You May Like, "Suggested For You" posts in News Feed, and People You May Know.⁷ And importantly, we also use AI to help keep our communities safe. For example, we use AI to help remove health-related misinformation,⁸ and since March 2020, we have removed over 18 million instances of COVID-19 misinformation.⁹

These are just a few examples of the many ways in which we use AI. We are invested in building new AI technologies, and more importantly, we are invested in building them responsibly. That is why Facebook created a dedicated, cross-disciplinary Responsible AI (RAI) team within its AI organisation. The RAI team works across the company's different product organisations, to build and test approaches to help ensure that our machine learning (ML) systems are designed and used responsibly.

The work of our RAI team is focused around five pillars that have been derived from consensus principles statements like those from the Commission and the OECD, to better guide our development of innovative new tools, frameworks, and processes to advance responsible AI.¹⁰ These five pillars are:

1. **Privacy & Security:** We are building new privacy-protective infrastructure to make it easier to enforce our privacy commitments across all our AI-driven products and services.
2. **Robustness & Safety:** We are building processes to help review and test some of our key AI systems prior to use in products and help ensure those systems behave safely and as intended even when they are subjected to malicious attack.
3. **Fairness & Inclusion:** We are building new tools and frameworks for testing the statistical fairness and fairness maturity of some of our key AI systems.
4. **Transparency & Control:** We are continuing to expand our industry-leading set of tools for explaining how our systems decide which feed content and ads to show you and are developing tools like "Model Cards"—which can provide simple, standardised documentation for many kinds of AI models.
5. **Accountability & Governance:** We have created a Responsible Innovation team that runs workshops with product teams during the development phase to raise awareness around individual and societal risks, and we are developing an AI-specific risk assessment framework.

The Responsible AI team is intended to help chart a course for the future of AI development that leads to a safer, fairer, and more prosperous society for all. We are also dedicating significant resources in building open research, knowledge, and tools that benefit AI researchers and entrepreneurs around the world. For example, we have invested in the AI research community in France and the EU through our Facebook AI Research (FAIR) hub in Paris. Our close collaboration with researchers on issues like computer vision, natural language processing, speech recognition, and more encourages greater innovation and helps bring private sector advances in research and technologies into academia.¹¹ We have also built open-source tools to advance AI research and development. In 2016, Facebook AI researchers, in collaboration with other AI leaders, created PyTorch, an open-source machine learning framework that embraces a philosophy of openness and collaborative research to advance state-of-the-art AI. PyTorch has become one of the default core development libraries used in the AI industry today. There are currently over 1,800 entities contributing to the PyTorch community—including institutions such as Caltech and companies pushing the boundaries of AI research, like OpenAI. We continue to make contributions to the PyTorch community through open-source projects like Captum, a library for model interpretability that helps researchers and developers better understand how AI models are making decisions. And these same open-source frameworks help Facebook do everything from improving augmented reality experiences to enabling Facebook AI Multimodal (FAIM), which allows us to identify harmful content across images, text, comments, and other elements holistically.¹²

We have learned important lessons about how to operationalise responsible AI principles and are contributing these insights to the global dialogue about AI governance. For example, we are collaborating with the OECD’s new AI Observatory project to study and disseminate emerging best practices that are in line with its 2019 AI principles. Similarly, through our recently launched Open Loop experimental governance program, we are conducting innovative “policy prototyping” approaches to co-create and test new potential AI policy frameworks with regulators and start-ups to ensure that future laws and regulatory instruments are both practical and evidence based. This consortium-driven effort has already launched projects in Europe, Asia, and Latin America, with more to come.¹³ And we are a founding partner in the Partnership on AI, the premier cross-industry, cross-civil-society multistakeholder forum for collaboratively developing AI best practices

III. The AI Act Must Allow for Refinement in Areas Where Foundational Questions Remain

Across many issues with AI and AI governance, clear or generally accepted standards, best practices, and frameworks do not yet exist. This uncertainty within the field and across industry has profound implications for the AI Act because in the absence of generally accepted standards and best practices, no one agrees on what it means for an AI system to be “fair”; no one knows the appropriate standards and measures for

concepts like robustness, security, accuracy, risk management, and risk assessment; no one has figured out how best to reconcile privacy with the need to measure bias through the collection of special categories of personal data; no one knows when a data set is complete; and no one really knows how these standards might change across different contexts, among other open questions. The AI Act should acknowledge this uncertainty and leave space for the development of standards and best practices that can fill these gaps.

A. Many of the AI Act's Requirements Are Not Feasible or Assume Best Practices or Standards That Do Not Yet Exist

Although the Commission's overall approach of focusing on high-risk AI systems and using self-assessments instead of third-party assessments strikes a good balance between regulation and innovation, many of the details of the Commission's approach threaten to limit AI innovation across Europe, and in some cases, may even prove to be counterproductive to the Commission's intended goals. In particular, many of the provisions could be interpreted to reach technologies that are either not high-risk AI systems, or are not even AI systems at all. Additionally, the Commission's approach to addressing bias in AI systems may not be technically feasible and may even unintentionally inhibit the development and use of certain fairness-promoting tools. This overbreadth could also make it difficult to develop and deploy even those lower-risk AI technologies that the Commission intends to regulate with a lighter touch.

1. The scope of many provisions reaches technologies that are not AI or not high-risk

The AI Act uses several key definitions that are overbroad and risk leading to overregulation of technologies that are not in fact high-risk AI systems or are not AI systems at all.

a) *Unclear definition of AI*

First, the proposed regulation uses too broad a definition of AI. Annex I of the proposed regulation defines AI to include "machine learning approaches" of various kinds, but also includes all "logic- and knowledge-based approaches," and any "statistical approach, Bayesian estimation, search and optimization methods."¹⁴ This definition of AI could potentially reach almost any modern software-based product because at some level all software is logic-based. The proposal's definition appears broad enough to cover everything from a simple digital wristwatch to the most complex AI systems.

We appreciate that the Commission is trying to craft a definition of AI that will be durable enough to encompass future technologies, but the current definition risks imposing new regulatory burdens on technologies that do not pose the sorts of risks the Commission intends to address through this regulation. We believe that the Commission should use a definition of AI that focuses on AI systems that learn and adapt over time because these are the capabilities that are at the core of AI, that make it different from other software applications, and that raise new and unique governance questions. The Commission's

definition of AI should distinguish between AI that learns to make decisions over time—including complex sets of decisions around human-level cognitive acts like driving a car, playing chess, or making judgments about someone’s job application—from software that produces outputs that are based on hard-coded, human-written rules.

To this end, we suggest the Commission consider a definition of AI that more closely aligns with globally accepted definitions that are already applied in industry. For example, the definition offered by the Expert Group on AI at the OECD captures the distinction between complex AI systems and general logic-based algorithms that we note above: “An AI system is a machine-based system that is capable of influencing the Environment by making recommendations, predictions or decisions for a given set of objectives. It does so by utilising machine and/or human-based inputs/data to: i) perceive real and/or virtual environments; ii) abstract such perceptions into models manually or automatically; and iii) use model interpretations to formulate options for outcomes.”¹⁵

b) Unclear prohibited practices

Second, the descriptions of the prohibited practices in Article 5 of the AI Act would benefit from further specificity. Combined with the significant penalties attached to violating those provisions, the unclear language could deter investment and development in many less risky or even beneficial technologies. For example, one of the prohibited practices includes “subliminal techniques.” But the AI Act does not define what “subliminal” means or what it means for such a technique to “materially distort a person’s behaviour in a manner that causes or is likely to cause that person or another person physical or psychological harm.”¹⁶ Prohibiting uses outright is the most blunt and drastic tool in the Commission’s toolbox, and the Commission is correct to limit its use to only the most dangerous and threatening uses of AI, but without clearly defining those uses, the Commission will end up preventing more AI applications than intended. For example, although Executive Vice President Vestager made clear that the Commission did not intend “subliminal techniques” to apply to targeted advertising,¹⁷ some observers have nonetheless argued that it could apply more broadly.¹⁸ By more carefully tailoring the language of the prohibited practices, the Commission can ensure that only the most dangerous uses of AI are prohibited.

c) Unclear scope of “remote biometric identification system”

Third, the scope of “remote biometric identification system” is also unclear. The use of AI to track and identify people through images and video without their consent or knowledge, and particularly when that is done by governments, poses risks to fundamental rights. At the same time, AI technologies using biometric data can have tremendous positive social benefits, where AI can be used to improve healthcare delivery and outcomes, monitor health and safety, enable people with disabilities to better navigate through the physical world, and more. The Commission must be careful that its definition of remote biometric identification does not limit the many positive uses of biometric data in AI. Two parts of the language would benefit from clarification. First, high-risk remote biometric identification systems are defined, in part, to be an “AI system for the purpose of identifying natural persons at a distance....”¹⁹ What it means to identify

someone “at a distance,” however, is not clear. This is particularly unclear because high-risk uses of remote biometric identification cover “post” (in addition to real-time) identification,²⁰ and it is difficult to imagine after-the-fact identification occurring in any way *other* than at a distance. For example, many people’s smartphones use AI to identify friends in photos, sometimes identifying people years after the photo was taken and hundreds or thousands of miles away from where the photo was taken or where the person is today. Therefore, such commonplace AI systems are identifying natural persons at a distance. Executive Vice President Vestager’s speech at the publication of the AI Act suggested that the intent was to cover mass surveillance where “many people are being screened simultaneously,”²¹ but the language should be clarified to reflect that intent.

It is also unclear what the Commission intended when it excluded from the definition of remote biometric identification where the “user of the AI system” has “prior knowledge ... whether the person will be present and can be identified.”²² For example, consumers might use their smartphone’s AI to find in their photos the faces of family and friends that they trained their device to recognise. In that example, it is unclear who the user of the AI system is. If the consumers are users, they arguably have “prior knowledge” whether the individuals in their contacts or their photo album can be identified by the device. But if the “user of the AI system” is the smartphone or software vendor that designed the AI system for the device, would they have prior knowledge? The Commission should clarify this language to not foreclose beneficial and common uses of AI to which people would be willing to consent, if given the appropriate opportunity.

d) Use of third-party conformity assessments and “safety components” as a proxy for high risk will be overbroad

Fourth, the Commission’s use of third-party conformity assessments in Union harmonisation legislation as a proxy for high risk could result in an overly broad application of “high risk” requirements. In addition to categorically classifying certain types of AI systems (those listed in Annex III) as “high risk,” Article 6 of the AI Act also classifies as high-risk AI systems that either would be a product, or would be a “safety component” of a product, that is covered by Union harmonisation legislation listed in Annex II, but only when the legislation requires the product to undergo third-party conformity assessments.²³

The issue, however, is that third-party conformity assessments may be mandated in Union sectoral legislation for reasons unrelated to the AI risk that the AI Act is addressing. This is particularly problematic with the inclusion of the Radio Equipment Directive 2014/53/EU (RED) in Annex II. Nearly all consumer electronics devices today make use of wireless communications, including cellular modems, Bluetooth, or WiFi, bringing them within scope of the Radio Equipment Directive. Under this Directive, a third-party assessment is required whenever there are no applicable harmonised standards for the product or where available harmonised standards can only be applied to the product in part. This situation is likely to arise with new and innovative products, where standards have often not yet been developed, or where conformity to applicable standards may need to be demonstrated through novel testing approaches. As a result, when a product

involving radio equipment that is required to undergo third-party conformity assessments under the Radio Equipment Directive includes AI, the AI Act will consider that system to be “high risk” even if the AI feature is unrelated to the reason for the required assessment.

Additionally, the definition and use of the term “safety component of a product” may actually discourage the use of AI to make systems safer when those systems are covered by legislation listed in Annex II. The AI Act appears to classify an AI system that serves *any* safety function, even if it is not satisfying a safety requirement established in the applicable Union harmonised legislation, as “high risk.” If an AI system that increases safety is classified as “high risk” by its mere presence in a product that is covered by legislation listed in Annex II—such as products including radio equipment—the AI Act may discourage the implementation of AI systems that improve safety by creating a significant regulatory barrier to bringing the feature to market.

e) Undifferentiated scope of “provider”

Finally, the AI Act uses an overly broad definition of “provider” that risks placing undifferentiated burdens on the diverse range of entities that participate in the AI development ecosystem, which could frustrate the Commission’s aims of supporting a healthy ecosystem of innovators, experimenters, contributors, and entrepreneurs. The AI Act defines a provider as any person or entity that “develops an AI system or that has an AI system developed with a view to placing it on the market or putting it into service under its own name or trademark, whether for payment or free of charge.”²⁴ This definition seems to assume that all AI systems are developed as a stand-alone product or service and then placed on the market. This assumption is wrong; the AI ecosystem is incredibly diverse, and there are many ways AI systems are developed and deployed. For example, the AI Act’s overly rigid provider and user taxonomy does not apply neatly to hybrid models such as where an AI user develops a system with a provider, a user develops AI in-house, or a user develops an AI system using more generalised models from an AI provider.

There is almost never a singular entity or person that develops an AI system. AI systems are the result of numerous entities building on top of others’ efforts. An AI system may start with open-source libraries, tools, and frameworks, created by hundreds or thousands of contributors (many not compensated) offering bits of code large and small.²⁵ Those open-source libraries, tools, and frameworks might then be combined with open data sets that themselves might be the work of hundreds of thousands of other people.²⁶ And the resulting model might then be shared under an open-source licence for others to build on. Who among all these contributors “develops an AI system?”

We urge the Commission to consider the range of developers, researchers, and innovators that make up the open-source community, which has been so crucial to advancing the state-of-the-art of AI development. Rather than a definition of “provider” that risks treating all contributions big and small to the same burdensome regulatory standards irrespective of their nature and role, we need a more nuanced taxonomy to

identify the relevant participants in the AI ecosystem and allocate the appropriate responsibilities and obligations to each one.

2. The proposal's attempt to promote fairness does not reflect the totality of causes of—and mitigations for—bias

Advancing AI systems that work fairly and equitably for everyone is incredibly important, and it is part of the reason why Facebook has invested in our Responsible AI team. We created the RAI team so that we could begin working toward an understanding of what it means for an AI system to be “fair,” to develop tools for measuring potential bias, to develop approaches for mitigating bias, and to create frameworks to help us balance competing equities and values—work that is still ongoing today.

The Commission's proposal attempts to ensure that high-risk AI performs well and is unbiased by mandating that data be “relevant, representative, free of errors, and complete.”²⁷ It is unclear whether any dataset can meet the Commission's proposed requirements. What we do know, however, is that “relevant, representative, free of errors, and complete” is not a feasible or practicable standard—and may actually make addressing bias harder.

Requiring that data be “relevant, representative, free of errors, and complete” has three significant problems: (a) bias does not come exclusively from data; (b) the standard is not technically achievable; and (c) it will limit the use of privacy and fairness enhancing technologies.

a) Bias does not come exclusively from data

The proposal appears mostly to address bias through the requirements that training data sets be error free and complete. This is problematic because it treats bad data as the primary source of bias. Bias in data sets could certainly be a problem and one that can be addressed in part through thoughtful data collection and curation practices. However, bias in AI systems does not come exclusively from biased or incomplete data. For example, data could be perfectly accurate and representative but reflect structural or societal biases. Or AI developers could make assumptions during the design of their systems as to how they expect the systems to be used, and as a result the system may exhibit bias. For example, an AI system that is built to try to predict criminality based on facial features will be biased from the outset due to faulty (and biased) assumptions in the design of the system, no matter what data is used. The point is that bias is not just a data issue but is one that must be addressed throughout the entire lifecycle of a product.²⁸

Instead of focusing solely on the data that is used to build an AI model (the “inputs”), the Commission should consider both inputs and outcomes holistically, recognising that error-filled inputs can still lead to fair outcomes, just as error-free inputs can still lead to unfair outcomes, depending on the overall design and operation of the system. A more holistic and systematic approach would be consistent with Executive Vice President Vestager's comments where she noted that the “proposed legal framework doesn't look

at AI technology itself. Instead, it looks at how AI is used, and what for.”²⁹ Particularly when it comes to bias, the Commission should consider the specific context in which an AI system is used, and whether the system’s design, inputs, and outputs are appropriate for that context, rather than focusing exclusively on data sets or any particular technology or methodology.

b) The requirements for training sets are unlikely to be technically feasible and could conflict with other important principles

Requiring data sets to be “free of errors” and “complete” is unlikely to be feasible. In practice no data set, particularly one that contains data collected about real world people and things, can ever be error free. This is true both for data sets used in supervised and unsupervised learning. In supervised learning, models are typically trained on pieces of data that are individually labelled by human annotators. This introduces two forms of potential error. First, even if there is clear agreement on what the label for any given element should be, human annotators will make mistakes, particularly when trying to efficiently label hundreds of thousands of items. But even if this kind of human error in annotation could be eliminated, we cannot eliminate the fact that in many cases there is no clearly “correct” label. Consider for example a model designed to detect hate speech; even the most highly trained expert annotators will find examples on which they will disagree. And yet we would not want to stop building tools to detect hate speech just because there may be ambiguity in the data.

The problems with this requirement are even more acute for unsupervised learning. Supervised learning is a labour-intensive process that is functionally limited by the capacity of human labellers to apply the (correct) labels to each element in the set. Although this works for smaller data sets, it cannot scale to large data sets or where accurate labels are non-existent or hard to find. For example, Facebook researchers have used unsupervised learning techniques to build speech recognition systems for languages such as Swahili and Tatar, which do not currently have high-quality speech recognition models available because they lack extensive collections of labelled training data.³⁰ Unsupervised learning is a way to train models using untagged data. Through numerous iterations, the model can discover clusters of data without human intervention.³¹ Training a language model, for instance, might use billions of examples of written language—far more than human annotators could provide—and that data might contain typos, poor grammar, and other errors. But given enough data, the model may still learn the correct spelling and grammatical structure of the language. If the AI Act foreclosed the use of imperfect data sets, this would de facto prohibit the use of unsupervised learning, which would be a significant impediment to AI innovation.

Unsupervised learning has significant advantages, and in particular it allows AI systems to be based on larger and more inclusive data sets than could be used if human annotation was required for all AI systems. For example, we might not have been able to build our Swahili and Tatar speech recognition systems as quickly (or at all) if we had to rely on

human annotation because accurate annotated data simply does not exist. Limiting the use of unsupervised learning could actually make AI less inclusive.

These issues aside, having error free data might not actually improve the functioning of AI systems. Good AI systems are designed to be robust even when they encounter unexpected or erroneous data. It is not clear that many or most AI systems would substantially improve if trained on error free data as opposed to data that is good enough. And if the costs of marginal improvements in the data set are substantial (which is often the case), the performance gains of the system are unlikely to be commensurate with the additional costs.

Similar problems arise with the requirement that data sets be “complete.” It is unclear what it would mean for a dataset to be “complete.” All datasets reflect choices—choices about time scale, population, geography, granularity, sample size, and so on. Emerging research on techniques such as directed sampling may eventually have a meaningful impact on improving models’ fairness in certain circumstances even where datasets are not complete, but this area of research is still nascent so best practices have yet to be established.³²

Even if it were technically feasible to have a complete data set, completeness is not always the best way to promote fair outcomes. For example, training a recommendation system on “complete” data that captures and reflects systemic bias in the real world would likely lead to a recommendation system that reflects that systemic bias. In that case, practitioners should not prioritise “complete” data sets when they would lead to less fair outputs but instead focus efforts on determining whether a recommendation system can be built that accounts for the underlying conditions that the data set reflects.

What’s more, requiring completeness may be at odds with other important principles that help safeguard fundamental rights, such as the data minimisation principle, which provides that data controllers should process only the personal data that is necessary to carry out a particular processing purpose. Requiring completeness where completeness is not necessary to fulfil a particular purpose would seem at odds with this principle (and potentially with the requirements of the GDPR).

c) The standard will limit the use of privacy and fairness enhancing tools

Over the last several years, a range of new tools have emerged that help protect privacy and advance fairness, and yet the “free of errors and complete” standard could actually make such tools impossible to use with AI systems. For example, privacy-preserving technologies like differential privacy are a way to share data and insights from data that minimises the risk of re-identification and enables the use of this data in the training of AI systems. Differential privacy tools take into account the sensitivity of the data set and add noise proportionally to ensure with high probability that no one can re-identify users.³³ “Noise” here is, in practice, errors that are purposely added to data sets such that individual data points cannot definitely be known. But because differential privacy injects

noise into the data, it could be argued that such techniques make the data less “free of errors” and their use in AI could be prohibited by the AI Act.

Similarly, the AI Act might stymie the use of approaches like the use of synthetic data, which is more privacy protective and helps ensure greater representativeness when real data is hard to collect. Synthetic data is artificially created data that can be used to train models to make them more robust and more accurate. This is particularly helpful as a mechanism for augmenting data sets with additional data that does not exist in the real world or would be difficult or cost-prohibitive to obtain. For example, Facebook recently developed the first multilingual machine translation model that can translate between any pair of 100 languages without relying on English data. In order to train the system to translate from Chinese to French, for example, it directly trained on Chinese to French data. In some cases, synthetic data was used to fill in gaps where there was insufficient direct translated data.³⁴ If we didn’t have synthetic data, the system would have been less accurate and less fair for speakers of less common languages.

Synthetic data could also involve creating combinations of data that are not as frequently observed, which can be important for mitigating bias in observed data. For example, if an employment data set includes data about female nurses, AI developers could create synthetic data about male nurses that would make the data set less skewed toward female nurses. Because this synthetic data would not directly and accurately reflect reality, it is arguably not “free of errors,” and the AI Act could prohibit this method of reducing bias.

Similarly, the standard may prevent the use of other techniques to debias data, a process by which data set providers intentionally alter their data to mitigate biases. For example, Spanish language texts might include more occurrences of “*la enfermera*” (feminine nurse) than “*el enfermero*” (masculine nurse). In this example, debiasing the data might involve changing some of the observed data to make the distribution of feminine and masculine terms more balanced. In both cases such changes might be considered introducing error into the data. If the AI Act effectively prohibited these practices, it would ultimately undermine the goal of reducing bias.

IV. We Need Collaboration to Develop Consensus Standards and Best Practices

A. Imposing Obligations Without Clear Standards Risks Serious Costs to Innovation

The costs of imposing the Act’s obligations in the absence of clear standards and best practices in key areas will be significant. It will create unintended consequences as innovators and entrepreneurs fear even low-risk AI systems might be swept into the overbroad and undefined categories.

AI is still an emerging technology where many of these standards and practices are still being developed. Building these norms, standards, and practices will require a concerted effort to coordinate between the many actors in the AI ecosystem. In fact, many of the AI Act's tools for updating the legislation (such as relying on implementing and delegating acts) do not provide for the involvement of all the relevant stakeholders. Only through the collaboration with—and engagement of—all stakeholders, can we strike the right balance between regulation and innovation, crafting standards that preserve fundamental rights whilst also being practical and achievable to ensure that crucial innovation is not hindered. It is this effort that the Commission should be seeking to support through the AI Act.

An example of this support for collaboration and the development of standards is the creation of regulatory “sandboxes,” which the proposal embraces.³⁵ While we are pleased to see the Commission's support for sandboxes, we believe the proposal could go even further in its support for a more holistic experimental governance framework to inform rulemaking processes and the development of standards and best practices. First, the AI Act's sandbox does not clearly identify the advantages for companies that participate in the sandbox, making it unclear why organisations would choose to join. Although the AI Act says that sandboxes should be a “safe space for experimentation,”³⁶ the sandbox provisions themselves do not offer any protections for companies that participate and actually note that sandboxes do not affect the supervisory powers of competent authorities,³⁷ and that participants remain liable for harm that results from experimentation in the sandbox.³⁸ And finally, the sandbox provisions offer no guidance as to how the insights from the sandboxes could feed back into the policymaking and regulatory process. Overall, this casts significant doubts as to whether the AI Act can sustain sandboxes that are attractive to organisations creating and developing new AI systems.

Second, sandboxes should be an opportunity for experimentation and co-creation for the purpose of enabling innovation both in technology development *and policymaking*. The Commission could provide support for the latter by more explicitly embracing “policy prototyping,” an experimentation-based approach for policy formulation and development that can provide a safe testing ground to try and learn how different policy approaches might play out when implemented in practice. Policy prototyping involves a variety of stakeholders that come together to co-create normative frameworks, including regulation and standards. By developing and testing regulatory frameworks in a collaborative fashion, this allows policymakers to see how their regulations can integrate with other co-regulatory tools such as corporate ethical frameworks, voluntary standards, certification programs, ethical codes of conduct, and best practices.

Facebook has been exploring how policy prototyping could support the development of more effective AI regulation through Open Loop.³⁹ Open Loop is a global strategic initiative that promotes and deploys experimental regulatory efforts in the field of new and emerging technologies. It supports the co-creation and testing of new governance frameworks through policy prototyping programs and enables the evaluation of existing

legal frameworks through regulatory sandbox exercises. Open Loop’s goal is to create a robust collaborative feedback loop (an “open loop”) of practical learnings between the people who make policy and those who must implement it. As part of Open Loop, policymakers work side-by-side with a vibrant community of tech companies, including Facebook, to build operational governance frameworks, and discuss regulatory best practices. Our initial work on policy prototyping through the Open Loop initiative has demonstrated the potential value in this kind of collaborative policymaking. As we describe in our first Open Loop report,⁴⁰ the process led to several concrete recommendations for improving self-assessments of AI. We would encourage the Commission to create more space in the regulation for efforts like these by explicitly acknowledging the role for policy prototyping and the development of standards that will provide much-needed guidance for providers subject to the Act.

In other words, what is needed from the Commission is not ironclad rules but mechanisms—such as sandboxes, policy prototyping and other consultative processes—for helping to answer some of the fundamental questions that need to be addressed. We have identified five key questions as a starting place.

1. What constitutes reasonable efforts to detect and mitigate bias in AI systems?

The Commission largely attempts to ensure that AI systems are fair by regulating the data that goes into those systems. As described above in Part III.2, we believe this is at best an insufficient approach, and at worst one that actually hampers efforts to make AI systems fairer. It is important to take a more comprehensive and practical approach to detecting and mitigating bias in AI systems, but generally accepted approaches for doing so in all circumstances simply do not yet exist. Therefore, the AI Act must leave space for the development of consensus methods for assessing, measuring, and comparing data and AI systems, as well as standards for reasonable mitigations. This will require the development of new frameworks, standards, and best practices.

There are many promising avenues that might someday soon become the foundation for these consensus approaches. For example, “model cards” could be a tool to help assess and compare AI models. Model cards are short documents that accompany AI models that provide standardised data about the model that make it easier for both experts and non-experts to understand how a model may perform under certain conditions or with certain kinds of data. Building on academic research in the area of model documentation,⁴¹ Facebook’s RAI team is collaborating with other product teams to develop and test our own standardised model card approach. For example, the Equity Team at Instagram used model cards as a way to monitor whether the AI systems used to limit the reach of misinformation or violent content were having a disproportionate impact on any one community.⁴² Similarly, data cards or data nutrition labels can help document the features and limitations of data sets, highlighting the potential risks and biases that may come from using a data set to train a particular model.⁴³

Model cards and data cards, however, are still being developed and refined as tools; researchers and industry are working to better understand the kinds of metrics that should be tracked and the best ways to highlight that data.

The AI Act sets out a rigid set of technical documentation that all high-risk AI systems must provide.⁴⁴ Instead of locking the EU into a narrow, inflexible set of documentation rules (the interplay of which with the GDPR is, in addition, uncertain), the AI Act should reconsider the level of prescriptiveness of the technical requirements and provide opportunities for stakeholders to come together to develop the standards and best practices for assessing, measuring, and comparing AI systems.

The appropriate documentation standards can help us better identify and measure bias in AI systems, but we also lack consensus standards for how to mitigate bias in AI systems once we've found it. Documentation alone will not determine whether an AI system is fair, nor will it tell us how to make systems *fairer*. Instead, we need to develop new standards and frameworks that help us determine reasonable efforts to take in the context of particular use cases. For example, while in certain cases addressing implementation errors in models can address certain types of unfair errors or performance disparities, in other cases addressing biased outcomes might be best achieved by including special categories of personal data in the model itself. Including such data, however, may actually contravene equality laws designed to prevent unfair treatment. Identifying appropriate approaches in cases like these will require a better understanding of the kinds of equitable outcomes that we as a society seek to achieve and the costs we are willing to bear to achieve them.

2. What are the appropriate standards and measures for robustness, accuracy, security, risk management and assessment?

The AI Act requires that “high-risk AI systems shall be designed and developed in such a way that they achieve, in the light of their intended purpose, an appropriate level of accuracy, robustness and cybersecurity, and perform consistently in those respects throughout their lifecycle.”⁴⁵ And yet there is no consensus about what “an appropriate level of accuracy, robustness, and cybersecurity” would mean. Is 87% accuracy “appropriate” or maybe 92%? Or perhaps the appropriate level changes with the specific context—but how do we know what the appropriate level is for critical infrastructure as opposed to a system used in vocational education? This is why it is important to leave space in the AI Act for the development of consensus best practices to establish what “appropriate” should mean.

Similar challenges arise with the AI Act’s requirements that high-risk AI providers implement risk management systems. The AI Act says that the risk management system should be a “continuous iterative process,”⁴⁶ but what is continuous? Every hour? Every week? Every year? Similarly, the provision says that mitigations should be used until the “overall residual risk of the high-risk AI systems is judged acceptable,”⁴⁷ or that “testing procedures shall be suitable.”⁴⁸ The use of language like “appropriate,” “acceptable,” or “suitable,” makes sense when regulating an area with extensive, existing best practices or

standards. That is not AI. We don't yet know what "appropriate," "acceptable," or "suitable," should be for high-risk AI systems (or any AI systems), and placing such terms into regulation—without ensuring they have widely accepted meanings—will create profound legal uncertainty for AI developers and stifle innovation. Instead, the Commission should encourage the multistakeholder processes necessary to develop best practices and standards that will help inform these concepts and defer to these processes to define the relevant terms in the AI Act.

Open Loop, our own experimental governance programme, works with companies in the EU to explore how something like risk assessment could work for AI systems.⁴⁹ Some of the lessons from our Open Loop project on risk assessment highlighted the need for policymakers and legislators to be very clear about the types of risk that they are attempting to address in any risk assessment requirement. Additionally, participants in our Open Loop research also emphasised the importance of developing clear procedures to help guide companies through the risk assessment process. Overall, this research reflects that much more work remains to be done to flesh out what something like risk assessment should look like and how to operationalise it.

3. How do we reconcile privacy with the need to measure bias using special categories of personal data?

It is impossible to correct a problem if it cannot be measured, detected, or monitored, and this presents a fundamental tension in addressing bias in AI systems. But detecting bias that affects vulnerable populations often requires collecting and storing special categories of personal data (e.g., race and ethnicity) so that system providers and maintainers can assess how the system is performing for certain categories of people.⁵⁰ But how can companies balance the privacy risk inherent in collecting and processing special categories of personal data for the purpose of measuring potential disparities, particularly if such data has not been self-provided or has not been provided for this specific purpose?

It is very helpful that the AI Act specifically allows practitioners to process special categories of personal data under the GDPR for purposes of measuring and mitigating bias in high-risk AI systems.⁵¹ However, this does not fully address this fundamental tension. First, the AI Act should extend this protection to all AI systems; measuring and mitigating bias is desirable in a wide range of AI systems, not just high-risk ones. For example, a system built to understand written language could learn undesired associations between offensive terms and certain populations of people, but if that system were not deemed high risk, developers may not be able to process the demographic data that would allow them to detect and mitigate that potential representational harm.

Second, allowing companies to process special categories of personal data does not address the question of how that data is collected in the first place—and whether for certain unobservable but protected characteristics, it is possible to collect at all.⁵² What systems can be put in place to help measure bias in AI systems without collecting or

processing special categories of personal data? And to the extent that such data must be collected, how can it be collected in the most privacy preserving manner? Here again, the AI Act should leave space so that diverse stakeholders can work together to collaboratively develop standards and best practices to ensure that the necessary data can be collected and used to better measure and mitigate bias.

4. How can we encourage the creation of more high-quality, open data sets?

As discussed above, it is unrealistic and infeasible to expect data sets to be “free of errors and complete.” But the Commission should be encouraging the creation of more high-quality and open data sets. Part of the reason why biased, low-quality data sets are used across the industry is that collecting diverse, representative, high-quality data is extremely difficult and expensive. Creating or commissioning such a dataset would likely be a barrier for a small start-up trying to create an innovative, new AI system.

Facebook understands, however, both the challenge and the potential of creating new, high-quality, open data sets. Recently, we published our Casual Conversations data set.⁵³ This is a data set consisting of 45,186 videos of paid actors having non-scripted conversations. By commissioning the data set, compensating participants, and using information that the participants provided directly to us, we have a data set that can be shared openly without risking exposing sensitive or private information. And most importantly, given the diversity of the data set, it is a critical tool to help the industry better test bias in their AI systems. Data sets like Casual Conversations can be extremely useful, but they are enormously expensive and time consuming to create. Rather than just dictating that data sets be “free of errors and complete,” the AI Act should leave space for the collaborative development of standards and best practices that can help establish measures for high-quality data sets. Additionally, the Commission should support the collection and maintenance of high-quality, publicly accessible data sets, so that all innovators and entrepreneurs can have access to free, high-quality data.

5. How can we encourage greater data sharing across industry and between the public and private sectors?

Throughout the AI Act, there are several provisions that would require AI providers to make available source code,⁵⁴ logs,⁵⁵ or training data.⁵⁶ We understand the importance of providing information necessary to enable regulators to ensure that AI systems are behaving appropriately and are in full compliance with the law. We believe that this kind of verification is critical to developing trustworthy AI systems. Such data sharing, however, must be done in a careful, privacy-preserving manner that does not create risks of adversarial attacks on AI systems.

Facebook is strongly committed to data transparency and data sharing where it has a positive impact for society. It is important to recognise, however, that full transparency into AI systems, and specifically the disclosure of source code, creates significant risks and may not actually help address the questions that regulators seek to answer. Full transparency may often not provide meaningful information to regulators. First, some of

our systems are the amalgamation of tens or hundreds of different models and algorithms and disclosing the source code of all of them would be unlikely to provide meaningful information to regulators. Second, the inner logic of our complex systems would entail the sharing of a vast amount of information that is irrelevant to helping regulators understand the outcomes. Additionally, the more widely sensitive operational data is shared, the greater the risks of accidental disclosure; such disclosure may allow bad actors to manipulate AI systems unethically or illegally for their own personal gain. Regulation should carefully tailor the data being shared to the purpose for which it is being shared. This is especially true in this case where the shared data may be sensitive and contain private information. The Commission should ensure that the AI Act provides sufficient privacy protections and outlines a process for careful vetting of recipients of shared data to minimise the risk of misuse and improper access.

B. Advancing the EU's Digital Single Market Strategy Requires Union-wide Consistency in Regulation and Enforcement

Building an effective digital single market across Europe in this area, which does not yet have industry-wide best practices or standards, requires special attention to both harmonisation and the expertise needed to effectively balance AI's risk and tremendous opportunity. Under the AI Act, Member States can designate existing authorities to oversee the implementation of the Act. Regulating AI, however, is likely to span multiple equities, from data protection to safety to economic, competition, and consumer issues. As such, the enforcement of the AI Act will likely require not only deep subject matter expertise but also navigation of the complex trade-offs that regulation will entail. As AI systems are built and regulated in Europe, common interpretations of risks, mitigations and the development of harmonised standards will also be essential to the growth of trusted, safe, fair, accountable, and transparent use cases across the region. In the absence of established and accepted best practices or standards on which to draw upon, mechanisms are especially needed to ensure that each supervisory authority is not left to create its own, divergent standards for enforcement. We encourage the Commission to consider enforcement and coordination mechanisms that enable a common regulatory environment across the digital single market and help foster sufficient expertise in and across the regulatory authorities.

V. Conclusion

Through our work at the forefront of building and governing AI technologies, we understand that these are still early days. We at Facebook—and we as a society—are at the beginning rather than the end of our responsible AI journey. In many cases we do not yet have answers to some of the most difficult questions that the use of AI poses. And in some cases, we do not yet have all the right questions.

Developing responsible, effective, practicable, and robust AI regulation requires first working to address these open questions, collaboratively building the frameworks, standards, and best practices that will be a necessary foundation and complement to future regulations. To that end, we look forward to working with the European

Commission and Members of European Parliament, and extend an invitation to join with us in our Open Loop policy prototyping programs and other collaborative endeavours that can help ensure that Europe's AI regulation sets a new global standard, unlocking a new era in responsible AI innovation for Europe.

We are excited about a future where AI is fair, transparent, accountable, and privacy-preserving, while also creating new economic opportunities, social interactions, and benefits for people around the world. For that reason, Facebook welcomes the European Commission's proposal and is eager to collaborate on the future of AI governance.

VI. Endnotes

- ¹ Zuckerberg, M. (2019) 'Four Ideas to Regulate the Internet', *About Facebook*, 30 March. Available at: <https://about.fb.com/news/2019/03/four-ideas-regulate-internet/> (Accessed: 20 February 2021).
- ² Fjeld, J. et al. (2020) 'Principled Artificial Intelligence: Mapping Consensus in Ethical and Rights-based Approaches to Principles for AI'. Available at: <https://dash.harvard.edu/handle/1/42160420> (Accessed: 4 June 2021).
- ³ High-Level Expert Group on AI, European Commission (2019) *Ethics guidelines for trustworthy AI, Shaping Europe's Digital Future*. Available at: <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai> (Accessed: 4 June 2021).
- ⁴ OECD (2019) *OECD Principles on Artificial Intelligence - Organisation for Economic Co-operation and Development*. Available at: <https://www.oecd.org/going-digital/ai/principles/?fbclid=IwAR0xJph4cqLWUvazLpUoSNK8iHrVXzimXWOr9IS-P8HF6a4JFvktakSMcLw> (Accessed: 17 February 2021).
- ⁵ Qiao, L. (2021) *PyTorch builds the future of AI and machine learning at Facebook*. Available at: <https://ai.facebook.com/blog/pytorch-builds-the-future-of-ai-and-machine-learning-at-facebook/> (Accessed: 4 June 2021).
- ⁶ Lada, A., Wang, M. and Yan, T. (2021) *How does News Feed predict what you want to see?*, Facebook Technology. Available at: <https://tech.fb.com/news-feed-ranking/> (Accessed: 16 February 2021).
- ⁷ Facebook, *What are recommendations on Facebook?* | Facebook Help Center. Available at: https://www.facebook.com/help/1257205004624246/?helpref=search&query=recommendation&s&search_session_id=5d5bca7781cbb1308fa641cf4818cf2a&sr=0 (Accessed: 15 March 2021).
- ⁸ Facebook AI (2020) *Using AI to detect COVID-19 misinformation and exploitative content*. Available at: <https://ai.facebook.com/blog/using-ai-to-detect-covid-19-misinformation-and-exploitative-content/> (Accessed: 16 February 2021).
- ⁹ Rosen, G. (2021) 'Moving Past the Finger Pointing', *About Facebook*, 17 July. Available at: <https://about.fb.com/news/2021/07/support-for-covid-19-vaccines-is-high-on-facebook-and-growing/> (Accessed: 19 July 2021).
- ¹⁰ Pesenti, J. (2021) 'Facebook's five pillars of Responsible AI', *Facebook AI*, 22 June. Available at: <https://ai.facebook.com/blog/facebook-five-pillars-of-responsible-ai/> (Accessed: 27 June 2021).
- ¹¹ 'Introducing Facebook AI Research Paris' (2015) *About Facebook*, 2 June. Available at: <https://about.fb.com/news/2015/06/introducing-facebook-ai-research-paris/> (Accessed: 28 June 2021).
- ¹² Qiao, L. (2021) *PyTorch builds the future of AI and machine learning at Facebook*. Available at: <https://ai.facebook.com/blog/pytorch-builds-the-future-of-ai-and-machine-learning-at-facebook/> (Accessed: 4 June 2021).
- ¹³ Andrade, N. and Kontschieder, V. (2021) *AI Impact Assessment: A Policy Prototyping Experiment*. Europe: Open Loop. Available at: <https://openloop.org/lets-unlock/>.
- ¹⁴ AI Act at Annex I.
- ¹⁵ OECD (2019) *OECD Principles on Artificial Intelligence - Organisation for Economic Co-operation and Development*. Available at: <https://www.oecd.org/going-digital/ai/principles/?fbclid=IwAR0xJph4cqLWUvazLpUoSNK8iHrVXzimXWOr9IS-P8HF6a4JFvktakSMcLw> (Accessed: 17 February 2021).
- ¹⁶ AI Act at Art. 5.

-
- ¹⁷ Vestager, M. (2021) ‘Speech by EVP Vestager at the press conference on AI’. Brussels, 21 April. Available at: https://ec.europa.eu/commission/presscorner/detail/en/speech_21_1866 (Accessed: 5 June 2021).
- ¹⁸ Douglas Heaven, W. (2021) ‘This has just become a big week for AI regulation’, *MIT Technology Review*, 21 April. Available at: <https://www.technologyreview.com/2021/04/21/1023254/ftc-eu-ai-regulation-bias-algorithms-civil-rights/> (Accessed: 31 July 2021).
- ¹⁹ AI Act at Art. 3(36).
- ²⁰ *Id.* at Annex III(1)(a).
- ²¹ Vestager, M. (2021) ‘Speech by EVP Vestager at the press conference on AI’. Brussels, 21 April. Available at: https://ec.europa.eu/commission/presscorner/detail/en/speech_21_1866 (Accessed: 5 June 2021).
- ²² AI Act at Art. 3(36).
- ²³ AI Act at Art. 6.
- ²⁴ *Id.* at Art.3(1)
- ²⁵ See, e.g., *pytorch/pytorch* (2021). pytorch. Available at: <https://github.com/pytorch/pytorch> (Accessed: 6 June 2021) (listing over 1,800 contributors to this open source library).
- ²⁶ See, e.g., ‘About ImageNet’, *ImageNet*. Available at: <https://image-net.org/about.php> (Accessed: 6 June 2021) (describing the collection of over 14 million images, labeled by Mechanical Turkers); *Open Images V6 - Description*. Available at: <https://storage.googleapis.com/openimages/web/factsfigures.html> (Accessed: 6 June 2021) (describing a collection of over 9 million images, computer labelled with some human verification); *CORD-19 / Semantic Scholar*. Available at: <https://www.semanticscholar.org/cord19> (Accessed: 6 June 2021) (describing a collection of 280,000 scholarly articles about COVID-19 research).
- ²⁷ AI Act at Art. 10(3).
- ²⁸ See Schwartz, R. *et al.* (2021) *A Proposal for Identifying and Managing Bias in Artificial Intelligence*. National Institute of Standards and Technology. doi: [10.6028/NIST.SP.1270-draft](https://doi.org/10.6028/NIST.SP.1270-draft), at 5-6 (discussing bias throughout the AI lifecycle).
- ²⁹ Vestager, M. (2021) ‘Speech by EVP Vestager at the press conference on AI’. Brussels, 21 April. Available at: https://ec.europa.eu/commission/presscorner/detail/en/speech_21_1866 (Accessed: 5 June 2021).
- ³⁰ Baevski, A. *et al.* (2021) *High-performance speech recognition with no supervision at all*. Available at: <https://ai.facebook.com/blog/wav2vec-unsupervised-speech-recognition-without-supervision/> (Accessed: 16 July 2021).
- ³¹ See IBM Cloud Education (2020) *What is Unsupervised Learning?* Available at: <https://www.ibm.com/cloud/learn/unsupervised-learning> (Accessed: 6 June 2021).
- ³² See, e.g., Rančić, S., Radovanović, S. and Delibašić, B. (2021) ‘Investigating Oversampling Techniques for Fair Machine Learning Models’, in Jayawickrama, U. *et al.* (eds) *Decision Support Systems XI: Decision Support Systems, Analytics and Technologies in Response to Global Crisis Management*. Cham: Springer International Publishing (Lecture Notes in Business Information Processing), pp. 110-123. doi: [10.1007/978-3-030-73976-8_9](https://doi.org/10.1007/978-3-030-73976-8_9).
- ³³ Herdağdelen, A. *et al.* (2020) ‘Protecting privacy in Facebook mobility data during the COVID-19 response’, *Facebook Research*, 3 June. Available at: <https://research.fb.com/blog/2020/06/protecting-privacy-in-facebook-mobility-data-during-the-covid-19-response/> (Accessed: 5 June 2021).

-
- ³⁴ Fan, A. (2020) 'Introducing the First AI Model That Translates 100 Languages Without Relying on English', *About Facebook*, 19 October. Available at: <https://about.fb.com/news/2020/10/first-multilingual-machine-translation-model/> (Accessed: 6 June 2021).
- ³⁵ AI Act at Art. 53.
- ³⁶ *Id.* at Recital 71.
- ³⁷ *Id.* at Art. 53(3).
- ³⁸ *Id.* at Art. 53(4).
- ³⁹ Open Loop, *A Global Experimental Governance Program*, Open Loop. Available at: <https://openloop.org/> (Accessed: 16 March 2021).
- ⁴⁰ Andrade, N. and Kontschieder, V. (2021) *AI Impact Assessment: A Policy Prototyping Experiment*. Europe: Open Loop. Available at: <https://openloop.org/lets-unlock/#europe-publication>.
- ⁴¹ Mitchell, M. *et al.* (2018) 'Model Cards for Model Reporting'. doi: [10.1145/3287560.3287596](https://doi.org/10.1145/3287560.3287596) ("Model cards are short documents accompanying trained machine learning models that provide benchmarked evaluation in a variety of conditions.")
- ⁴² Brooks, R., Houston, S. and Wadsworth, C. (2021) *Equity updates – what we've been up to and what comes next*, *Instagram Blog*. Available at: <https://about.instagram.com/blog/announcements/equity-updates-what-weve-been-up-to-and-what-comes-next> (Accessed: 15 July 2021).
- ⁴³ See Bender, E. M. and Friedman, B. (2018) 'Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science', *Transactions of the Association for Computational Linguistics*, 6, pp. 587–604. doi: [10.1162/tac1_a_00041](https://doi.org/10.1162/tac1_a_00041).
- ⁴⁴ See AI Act, Annex IV.
- ⁴⁵ *Id.* at Art. 15(1).
- ⁴⁶ *Id.* at Art. 9(2).
- ⁴⁷ *Id.* at Art. 9(4).
- ⁴⁸ *Id.* at Art. 9(6).
- ⁴⁹ Andrade, N. and Kontschieder, V., (2021) *AI Impact Assessment: A Policy Prototyping Experiment*, available at: https://d32j3j47emgb6f.cloudfront.net/wp-content/uploads/2021/01/AI_Impact_Assessment_A_Policy_Prototyping_Experiment.pdf
- ⁵⁰ Bogen, M., Rieke, A. and Ahmed, S. (2020) 'Awareness in Practice: Tensions in Access to Sensitive Attribute Data for Antidiscrimination', *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 492–500. doi: [10.1145/3351095.3372877](https://doi.org/10.1145/3351095.3372877).
- ⁵¹ AI Act at Art. 10(5)
- ⁵² Tomasev, N. *et al.* (2021) 'Fairness for Unobserved Characteristics: Insights from Technological Impacts on Queer Communities', *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 254–265. doi: [10.1145/3461702.3462540](https://doi.org/10.1145/3461702.3462540).
- ⁵³ Hazirbas, C. *et al.* (2021) *Shedding light on fairness in AI with a new data set*. Available at: <https://ai.facebook.com/blog/shedding-light-on-fairness-in-ai-with-a-new-data-set/> (Accessed: 6 June 2021).
- ⁵⁴ See, e.g., AI Act at Art. 64(2).
- ⁵⁵ *Id.* at Art. 23.
- ⁵⁶ *Id.* at Art. 64(1).