



THE UNIVERSITY OF
MELBOURNE

Transparency: *Decisions & Processes*

Marc Cheong

—
School of Computing and Information Systems
Centre for AI & Digital Ethics
The University of Melbourne
marc.cheong [at] unimelb.edu.au





Learning Outcomes

1. Distinguish between transparency and explainability, closely-related concepts in AI ethics.
2. Understand how automated decision-making (ADM) systems require transparency at every stage.
3. Understand how tech companies approach the issue of transparency, as well as concerns that are raised by stakeholders, in two areas where AI systems are deployed: social media ads and criminal justice.
4. Understand how big data research can - either positively or negatively - affect their data subjects, and why transparency is important when conducting such studies.



Related Reading

This module has two readings corresponding to the two broad themes within (plus an optional study).

Recap: screenshot below 😊

1. [Experiments in Social Media ↗](#)

Toby Walsh. *AI Magazine*, 40(4), 74-77. 2019. Archived copy by the author in arXiv [cs.CY].

This is an interesting piece which highlights the dangers how studies/experiments on social media can contribute to "challenges[,] as even small effects when multiplied by a large population can have a significant impact". Experimental "interventions increased turnout by about 340,000 additional votes ... around 0.5% of the total number of votes cast" (Walsh, 2019) in a Facebook experiment to encourage voting in the 2010 US Elections. The author highlights the issues of (non)transparency in AI research especially on such a large scale.

2. [Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead ↗](#)

Cynthia Rudin. *Nature Machine Intelligence*, 1, 206–215. 2019. Archived copy by the author in arXiv [stat.ML].

The second reading for this week focuses on the pervasiveness of "black box machine learning models" (Rudin, 2019) and how their implementation -- in, say, decision making for criminal justice -- can be mired by a "lack of transparency and accountability of predictive models ... [with] severe consequences" (Rudin, 2019). We need to look at the processes and the decision making aspects behind the development and deployments of these systems. From a computer scientist's lens, we should ask questions like 'should we even use it on people?' rather than 'can we optimise it?'

Additional readings.

3

If you find this module interesting, you might want to check out the brief history of the [Cambridge Analytica controversy on Wikipedia ↗](#) here.



Outline

1. Transparency: what's it all about.
2. Automated decision-making (ADM) systems: transparency at every stage?
3. Current issues in transparency and stakeholders' concerns,
 - I. social media advertising systems
 - II. criminal justice AI systems
4. Data science research and transparency: the cases of Cambridge Analytica and Covid19 Trends.



THE UNIVERSITY OF
MELBOURNE



Transparency: What's it all about?

Transparency? (1/2)



ENGLISH ▾



Privacy expert argues “algorithmic transparency” is crucial for online freedoms at UNESCO knowledge café

2 min

As more decisions become automated and processed by algorithms, these processes become more opaque and less accountable, with risks of secret profiling and illegal discrimination. For Rotenberg, “at the core of modern privacy law is a single goal: to make transparent, the automated decisions that impact our lives.” He sees “algorithmic transparency”, the principle that data processes which impact individuals be made public, as the next stage in the development of transparency law, internet law and privacy law. The lack of algorithmic transparency in the current internet ecosystem poses a crucial challenge to defending fundamental human rights online, ranging from privacy and freedom of expression to security. In addition to algorithmic transparency, Rotenberg pointed to other emerging issues which need to be examined, notably the increasing access to drones and robots and the need for their registration.

Source: UNESCO news release, quoting Marc Rotenberg: <https://en.unesco.org/news/privacy-expert-argues-algorithmic-transparency-crucial-online-freedoms-unesco-knowledge-cafe>

Transparency? (2/2)

Institutional transparency and public values

There are many dimensions to algorithmic ‘transparency’, but in the context of institutional actors, it requires clarity in the procurement, implementation and technical mechanisms associated with automated decision-making systems. This type of transparency is useful for keeping track of the impacts of decision systems over time, and achieving some public disclosure on their purpose, reach, policies, and techniques.

Sources:

Jake Goldenfein, 'Algorithmic Transparency and Decision-Making Accountability: Thoughts for buying machine learning algorithms' in Office of the Victorian Information Commissioner (ed), Closer to the Machine: Technical, Social, and Legal aspects of AI (2019).
https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3445873



A note on nomenclature (and focus)

Some academics use the term transparency to describe the “inner workings” of algorithms (e.g. Loi et al, 2020: <https://link.springer.com/article/10.1007/s10676-020-09564-w>)

To clarify, the focus of this module is on the overarching processes and decision-making of an algorithmic system.

We include the technical considerations about interpretability of the algorithms themselves, counterfactual analysis, etc under the banner of explainability.

(Our very own Tim Miller is an expert in this field).



Example: Hypothetical Thought Experiment. Deep Learning and Your Grades!

What if, for this unit, we decide your final grade based on a state-of-the art deep learning-based estimator/predictor that will come up with a final grade based on 1,000 factors - ranging from how many Youtube videos you watch about ethics, to your activity in weekly discussion activities, to your typing speed when asked to comment on discussion boards, to your on-Zoom reactions to Simon and the tutors when they taught you the basics of ethics, to how much you laugh at Marc's internet memes in the modules, etc. The final mark predictor is so advanced, it has been audited by NASA, Google, and by 10 Nobel Prize winners! However, since some of the tensor-based algorithms used is proprietary to NVidia, who sponsored the array of Geforce RTX3090s used for deep learning, they can't be revealed in public. Also, the decision made by the predictor is final.

Reflection.

What if, for this unit, we decide your final grade based on a state-of-the art deep learning-based estimator/predictor that will come up with a final grade based on 1,000 factors - ranging from how many Youtube videos you watch about ethics, to your activity in weekly discussion activities, to your typing speed when asked to comment on discussion boards, to your on-Zoom reactions to Simon and the tutors when they taught you the basics of ethics, to how much you laugh at Marc's internet memes in the modules, etc. The final mark predictor is so advanced, it has been audited by NASA, Google, and by 10 Nobel Prize winners! However, since some of the tensor-based algorithms used is proprietary to NVidia, who sponsored the array of Geforce RTX3090s used for deep learning, they can't be revealed in public. Also, the decision made by the predictor is final.

Process: who governs the selection of ML models/training data?

What vendors are given preference – e.g. why Google (not AWS)?

e.g why Tensorflow and Nvidia?

Any conflicts of interest? Any feedback loops?

Why is the whole system shrouded in secrecy?

Who audited it? Can I see the source code/design schematics/rationales?

Did I even sign up for this?

Decisions: can we challenge them? Can I take this to court?

Who sanctioned this to be official?

Is this another *Cambridge Analytica*?





THE UNIVERSITY OF
MELBOURNE



Automated decision-making (ADM) systems: *Transparency at every stage?*

Disclaimer: I am not a lawyer



The information provided in this mini-lecture is summarized from various sources to explain how transparency is a requirement not only for the algorithms, but also the contexts surrounding their implementation.

This lecture won't make you an expert in administrative decision-making ☺

Automated Assistance in Administrative Decision Making

The report contains best practice principles for the development and operation of expert computer systems used to make or assist in the making of administrative decisions. The Council believes the principles will ensure that decisions made using expert systems are consistent with existing administrative law values.

Commonwealth Ombudsman [Automated decision-making better practice guide](#)

Source: Administrative Review Council, *Automated Assistance in Administrative Decision Making: Report to the Attorney General* (Report No 46, November 2004) ('2004 Report')

<https://www.ag.gov.au/legal-system/publications/report-46-automated-assistance-administrative-decision-making-2004>

The 2004 Report was ahead of its time (emphases below are mine)

P27: "**Safeguards built into the system are only asking relevant questions, telling customers why questions are being asked (which makes the decision-making process more transparent)** and recording and explaining to a customer the reason for a decision"

P43: "Expert systems' ability to provide an **audit trail of the administrative decision-making processes they are involved in** is important to the administrative law values of transparency, fairness and efficiency."

P45: "A good system of internal review is one which **is transparent in process and affords a quick, inexpensive and independent review of decisions**. Such a system is beneficial both to applicants and agencies". (citing Administrative Review Council 2000)

GDPR?

4. Transparency and accountability in the General Data Protection Regulation

A number of provisions in the GDPR seek to promote a high degree of transparency in the processing of personal data [6]. In general these provisions require data controllers to provide data subjects with information about the processing of their personal data and to do so in a concise, transparent, intelligible and easily accessible form, using clear and plain language.

Where personal data are obtained from the data subject, Article 13(2)(f) requires data controllers to provide data subjects with information about ‘the existence of automated decision-making, including profiling ... and meaningful information about the logic involved, as well as the significance and the envisaged consequences of such processing for the data subject.’ The purpose of such information provision is said to be ‘to ensure fair and transparent processing’.

Sources:

Christina Blacklaws, ‘Algorithms: transparency and accountability’, Phil. Trans. R. Soc. A.3762017035120170351. (2018).
<https://royalsocietypublishing.org/doi/10.1098/rsta.2017.0351>

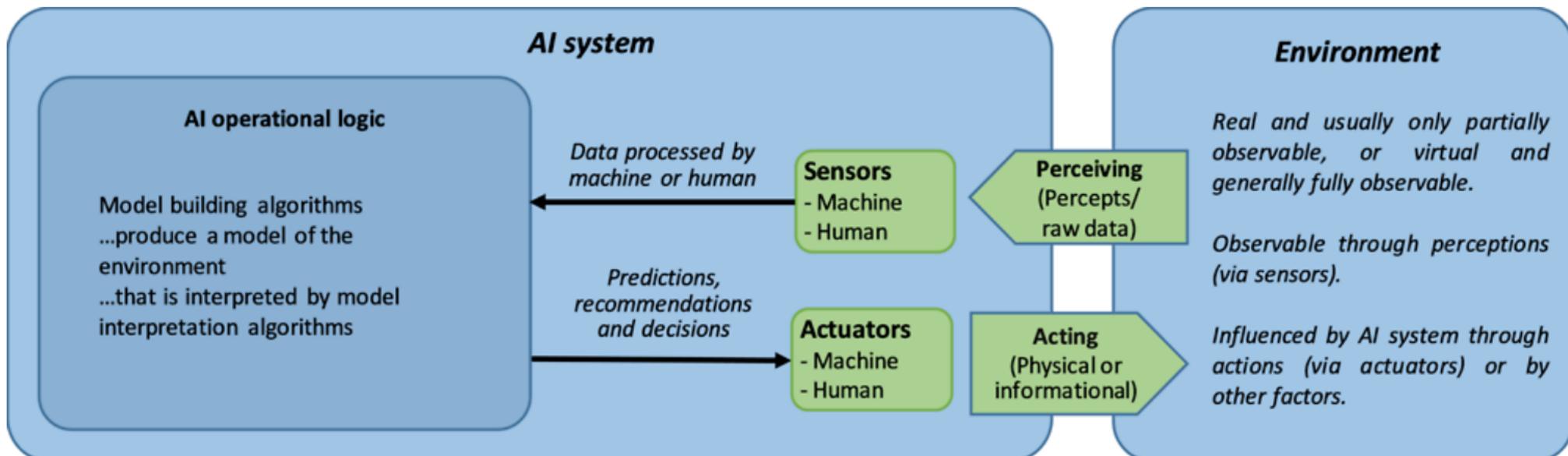
High-level view of AI systems (OECD, 2019)

Source: OECD (2019), *Artificial Intelligence in Society*, OECD Publishing, Paris, <https://doi.org/10.1787/eedfee77-en>.

Figure: <https://www.oecd-ilibrary.org/sites/8b303b6f-en/index.html?itemId=/content/component/8b303b6f-en#figure-d1e976>



Figure 1.3. A high-level conceptual view of an AI system



Source: As defined and approved by AIGO in February 2019.

Elements for ML systems (Google Inc, from Scully et al 2015)

Source: <https://cloud.google.com/solutions/machine-learning/mlops-continuous-delivery-and-automation-pipelines-in-machine-learning> - adapted from Scully et al (2015) <https://papers.nips.cc/paper/2015/file/86df7dcfd896fcf2674f757a2463eba-Paper.pdf>

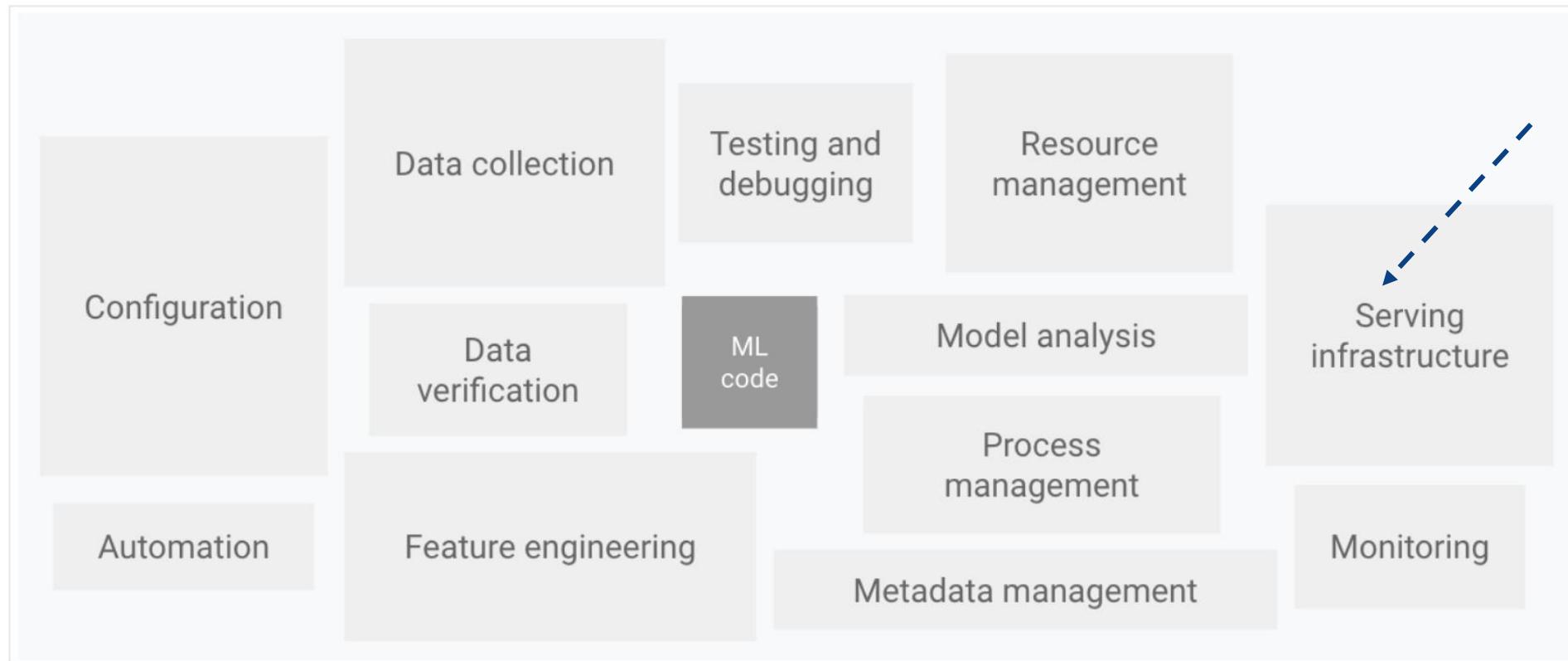


Figure 1. Elements for ML systems. Adapted from [Hidden Technical Debt in Machine Learning Systems](#).

Reflection.

From our examples, the code ('algorithm') is only a small part of it!

Transparency is required in planning, implementation, auditing...

... concerns the data, design, <actual ML stuff here>, testing, deployment

... also consider legal aspects & philosophical concepts in the broader sense: incl. fairness, recourse...





THE UNIVERSITY OF
MELBOURNE

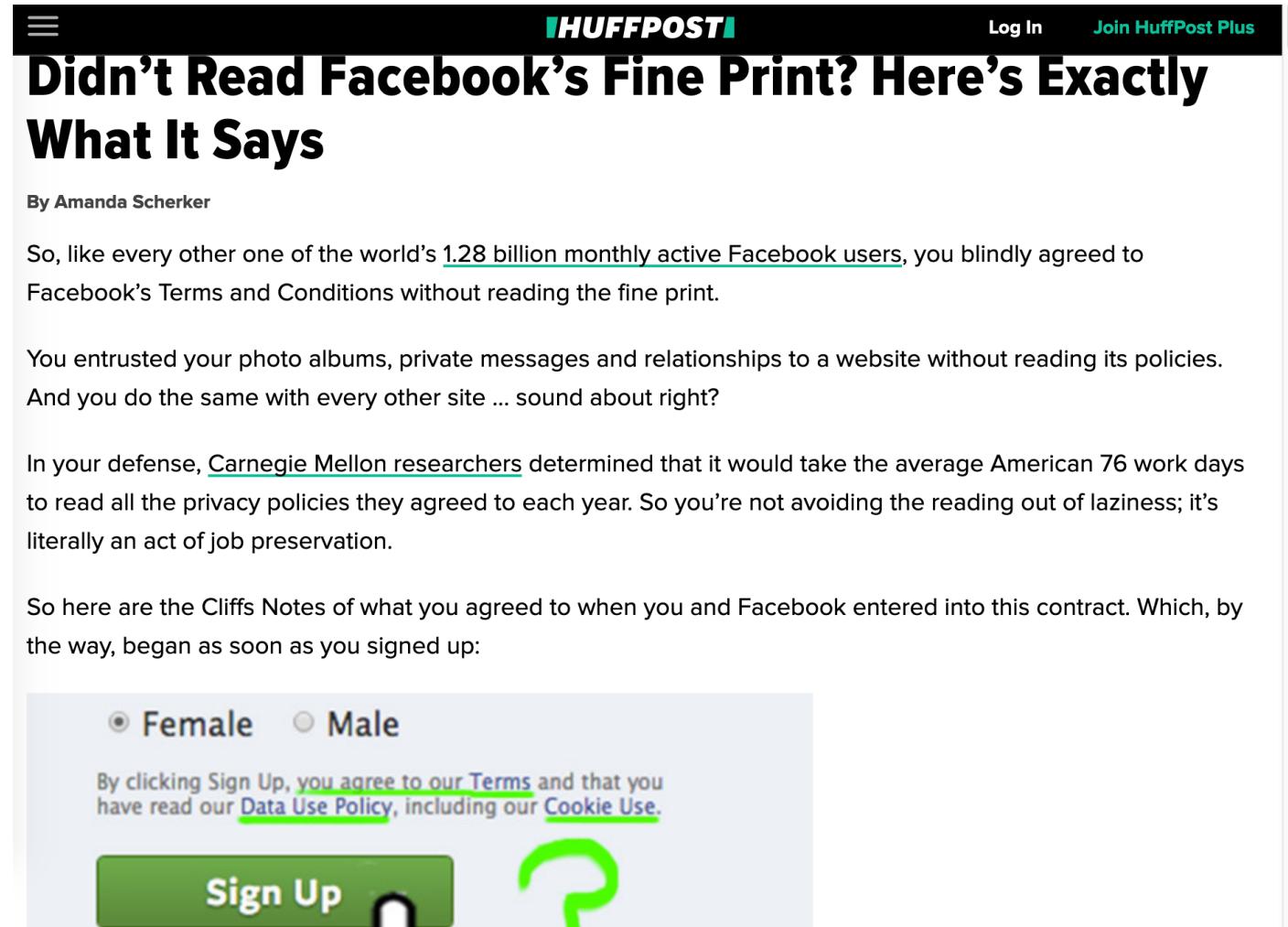


Current issues in transparency: *Social media advertising systems*

Consider Facebook's group of products

If we use Instagram,
Facebook, WhatsApp,
how many of us actually
read the TOS?

Source: Scherker (2017)
https://www.huffpost.com/entry/facebook-terms-condition_n_5551965



The screenshot shows a news article from HuffPost titled "Didn't Read Facebook's Fine Print? Here's Exactly What It Says" by Amanda Scherker. The article discusses how users of various platforms like Instagram, Facebook, and WhatsApp likely agreed to their terms of service without reading them. It引用了Carnegie Mellon研究人员的发现，指出平均美国人需要76个工作日才能阅读所有隐私政策。文章还展示了Facebook注册界面的一个片段，显示了性别选择（女性或男性）、同意条款和数据使用政策的勾选框，以及“Sign Up”按钮和一个带有问号的图标。

HUFFPOST

Log In Join HuffPost Plus

Didn't Read Facebook's Fine Print? Here's Exactly What It Says

By Amanda Scherker

So, like every other one of the world's 1.28 billion monthly active Facebook users, you blindly agreed to Facebook's Terms and Conditions without reading the fine print.

You entrusted your photo albums, private messages and relationships to a website without reading its policies. And you do the same with every other site ... sound about right?

In your defense, Carnegie Mellon researchers determined that it would take the average American 76 work days to read all the privacy policies they agreed to each year. So you're not avoiding the reading out of laziness; it's literally an act of job preservation.

So here are the Cliffs Notes of what you agreed to when you and Facebook entered into this contract. Which, by the way, began as soon as you signed up:

Female Male

By clicking Sign Up, you agree to our [Terms](#) and that you have read our [Data Use Policy](#), including our [Cookie Use](#).

[Sign Up](#) 



Consider Facebook's group of products

Roughly 9-10 pages in A4 printed (as of time of this lecture).

Summarised by TOSDR
(<https://tosdr.org/en/service/182>) - even the key points take up ~2 A4 pages;
refer image →

This doesn't even include the additional policies: e.g. 'Commercial Standards', 'Advertising Policies, ... (~12 more links)

Tracking pixels used in service-to-user communication
Facebook uses cookies
Facebook uses your data for many purposes
App required for this service requires broad device permissions
Your identity is used in ads that are shown to other users
personal data is given to third parties
The service uses your personal data for advertising
Your biometric data is collected
Users are not allowed to use pseudonyms, as trust and transparency between users regarding their identities is relevant to the service.
This service tracks you on other websites
The service can read your private messages
Facebook stores your data whether you have an account or not.
This service gathers information about you through third parties
The service informs users that its privacy policy does not apply to third party websites
There is a date of the last update of the terms
You can choose with whom you share content
Invalidity of any portion of the Terms of Service does not entail invalidity of its remainder
Usernames can be rejected for any reason
Users who have been permanently banned from this service are not allowed to re-register under a new account
This service may collect, use, and share location data
The service will not allow third parties to access your personal information without a legal basis
Failure to enforce any provision of the Terms of Service does not constitute a waiver of such provision



Let's try Twitter?

Roughly 39 pages in A4 printed!!! (as of time of this lecture).

https://cdn.cms-twdigitalassets.com/content/dam/legal-twitter/site-assets/tos-oct-14th-2020/Twitter_User_Agreement_EN.pdf

Summarised by TOSDR
(<https://tosdr.org/en/service/195>) - even the key points take up ~2 A4 pages;
refer image →

The screenshot shows the TOSDR service page for Twitter. At the top, there is a large blue Twitter logo with the word "Twitter" below it. To the right of the logo is a red box containing the word "Grade E". Below the logo, there is a URL "https://shields.tosdr.org/en_195.svg". On the left side of the page, there is a small red box with the text "tosdr.org/#twitter Privacy Grade E". The main content area is a vertical list of statements about Twitter's terms of service, each followed by two small circular icons.

Statement	Icons
Third party cookies	info link icon, link icon
You can retrieve an archive of your data	info link icon, link icon
This service is only available to users of a certain age	info link icon, link icon
You are responsible for maintaining the security of your account and for the activities on your account	info link icon, link icon
You have the right to leave this service at any time	info link icon, link icon
Your data may be processed and stored anywhere in the world	info link icon, link icon
They may stop providing the service at any time	info link icon, link icon
This service prohibits users from attempting to gain unauthorized access to other computer systems	info link icon, link icon
Users should revisit the terms periodically, although in case of material changes, the service will notify	info link icon, link icon
Failure to enforce any provision of the Terms of Service does not constitute a waiver of such provision	info link icon, link icon
This service can license user content to third parties	info link icon, link icon
Invalidity of any portion of the Terms of Service does not entail invalidity of its remainder	info link icon, link icon
The service uses your personal data for advertising	info link icon, link icon
The service can delete specific content without prior notice and without a reason	info link icon, link icon
The service is provided 'as is' and to be used at the users' sole risk	info link icon, link icon
This service reserves the right to disclose your personal information without notifying you	info link icon, link icon
This service provides a way for you to export your data	info link icon, link icon
There is a date of the last update of the terms	info link icon, link icon
The service allows you to use pseudonyms	info link icon, link icon
If you are the target of a copyright claim, your content may be removed	info link icon, link icon

Third party cookies	info link icon, link icon
You can retrieve an archive of your data	info link icon, link icon
This service is only available to users of a certain age	info link icon, link icon
You are responsible for maintaining the security of your account and for the activities on your account	info link icon, link icon
You have the right to leave this service at any time	info link icon, link icon
Your data may be processed and stored anywhere in the world	info link icon, link icon
They may stop providing the service at any time	info link icon, link icon
This service prohibits users from attempting to gain unauthorized access to other computer systems	info link icon, link icon
Users should revisit the terms periodically, although in case of material changes, the service will notify	info link icon, link icon
Failure to enforce any provision of the Terms of Service does not constitute a waiver of such provision	info link icon, link icon
This service can license user content to third parties	info link icon, link icon
Invalidity of any portion of the Terms of Service does not entail invalidity of its remainder	info link icon, link icon
The service uses your personal data for advertising	info link icon, link icon
The service can delete specific content without prior notice and without a reason	info link icon, link icon
The service is provided 'as is' and to be used at the users' sole risk	info link icon, link icon
This service reserves the right to disclose your personal information without notifying you	info link icon, link icon
This service provides a way for you to export your data	info link icon, link icon
There is a date of the last update of the terms	info link icon, link icon
The service allows you to use pseudonyms	info link icon, link icon
If you are the target of a copyright claim, your content may be removed	info link icon, link icon



Consider ads: is this ‘transparent’?

Let's say I've been shown an ad (sponsored post) on Facebook. As a consumer, I want to know WHY.

(Image sources: Facebook).

I go to the TOS, then find 'Ads', then find 'Data Policy'...

Which takes me to another page...

Which scares me.

Let's try another method – go to the ad, and click on the menu.

Device Information

As described below, we collect information from and about the computers, phones, connected TVs and other web-connected devices you use that integrate with our Products, and we combine this information across different devices you use. For example, we use information collected about your use of our Products on your phone to better personalize the content (including ads) or features you see when you use our Products on another device, such as your laptop or tablet, or to measure whether you took an action in response to an ad we showed you on your phone on a different device.

Information we obtain from these devices includes:

- **Device attributes:** information such as the operating system, hardware and software versions, battery level, signal strength, available storage space, browser type, app and file names and types, and plugins.
- **Device operations:** information about operations and behaviors performed on the device, such as whether a window is foregrounded or backgrounded, or mouse movements (which can help distinguish humans from bots).
- **Identifiers:** unique identifiers, device IDs, and other identifiers, such as from games, apps or accounts you use, and Family Device IDs (or other identifiers unique to Facebook Company Products associated with the same device or account).
- **Device signals:** Bluetooth signals, and information about nearby Wi-Fi access points, beacons, and cell towers.
- **Data from device settings:** information you allow us to receive through device settings you turn on, such as access to your GPS location, camera or photos.

Consider ads: is this ‘transparent’?

Let's say I've been shown an ad (sponsored post) on Facebook. As a consumer, I want to know WHY.
(Image sources: Facebook)

Why You're Seeing This Ad

Only you can see this

You're seeing this ad because your information matches [REDACTED] advertising requests. There could also be more factors not listed here. [Learn More](#)

[REDACTED] is trying to reach people, ages 18 and older. >

[REDACTED] is trying to reach people whose primary location is Australia. >

What You Can Do

Hide all ads from this advertiser

You won't see [REDACTED] ads

Make changes to your ad preferences

Adjust settings to personalize your ads >

Location information:

We may use your location information to show you ads from advertisers trying to reach people in or near a specific place. We get this information from things such as:

- Where you connect to the internet.
- Where you use your phone.
- Your location from your Facebook and Instagram profile.



Consider ads: is this ‘transparent’?

Case study thanks to Michael Geers,
(Max Planck Institute for Human Development, Germany)

Let's say I've been shown an ad (sponsored post) on Facebook.
As a consumer, I want to know WHY.

FB gives me some reasons.

... but also “more factors not listed”.

... and directs me to a huge page with many explanations,
but all “mays” (we may, advertisers may...)

... and has a long TOS explaining the many ways.



Question: Is this transparent enough?
What should/can social media companies do?



Consider ads: is this ‘transparent’?

Case study thanks to Michael Geers,
(Max Planck Institute for Human Development, Germany)

A simple self-reflection intervention boosts the detection of microtargeted advertising

AUTHORS

Philipp Lorenz-Spreen, Michael Geers, Thorsten Pachur, Ralph Hertwig, Stephan Lewandowsky, Stefan Herzog

<https://psyarxiv.com/ea28z/> - Lorenz-Spreen et al (2021)

P4-5:

“At present, the platforms’ transparency measures offer “nominal transparency”, with no real regard for whether people actually can easily access, read and gain insight into the information held about them and whether this transparency in name foster users’ autonomy.

“Aiming for effective transparency—which demonstrably enables users to understand what platforms do with their data and what users’ choices imply, and to then translate this knowledge into measurable behaviour—is an important step towards more acceptable business practices and towards regaining some of the lost autonomy for users (e.g., by prompting people to adjust their privacy settings; Parra-Arnau et al., 2017)



THE UNIVERSITY OF
MELBOURNE



Current issues in transparency: *Criminal justice AI systems*



Criminal Justice and AI systems

You might have encountered the following examples in your exploration of AI ethics.
(Images are from Wikipedia)

COMPAS (software)

From Wikipedia, the free encyclopedia

Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) is a [case management](#) and [decision support tool](#) developed and owned by Northpointe (now [Equivant](#)) used by [U.S. courts](#) to assess the likelihood of a [defendant](#) becoming a [recidivist](#).^{[1][2]}

COMPAS has been used by the U.S. states of New York, Wisconsin, California, Florida's [Broward County](#), and other jurisdictions.^[3]

PredPol

From Wikipedia, the free encyclopedia

PredPol, Inc is a [predictive policing](#) company that attempts to predict property crimes using [predictive analytics](#).

PredPol is also the name of the software the company produces. PredPol began as a project of the [Los Angeles Police Department](#) (LAPD) and [UCLA](#) professor Jeff Brantingham. PredPol has produced a patented algorithm, which is based on a model used to predict earthquake [aftershocks](#).

As of 2020, PredPol's algorithm is the most commonly used predictive policing algorithm in the U.S.^{[1][2]} Police departments that use PredPol are given printouts of jurisdiction maps that denote areas where crime has been predicted to occur throughout the day.^[3] The [Los Angeles Times](#) reported that officers are expected to patrol these areas during their shifts, as the system tracks their movements via the GPS in their patrol cars.^[4] Scholar Ruha Benjamin called PredPol a "crime production algorithm," as police officers then more heavily patrol these predicted crime zones, expecting to see crime, which leads to a self-fulfilling prophecy.^[1]

PredPol

Type	Private
Headquarters	Santa Cruz
Products	Predictive analytics
Website	www.predpol.com



Reading: Rudin (2019) on COMPAS

The focus of this case study is not about the technical explainability for the underlying algorithms, etc.

... but about the transparency of the processes and decisions involved.

Take COMPAS – and its **decision making** assumptions, in practice (Rudin 2019)

Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use
Interpretable Models Instead

Cynthia Rudin
Duke University

But if the model is a black box, it is very difficult to manually calibrate how much this additional information should raise or lower the estimated risk. This issue arises constantly; for instance, the proprietary COMPAS model used in the U.S. Justice System for recidivism risk prediction does not depend on the seriousness of the current crime [27, 29]. Instead, the judge is instructed to somehow manually combine current crime with COMPAS. Actually, it is possible that many judges do not know this fact. If the model were transparent, the judge could see directly that the seriousness of the current crime is not being considered in the risk assessment.



Reading: Rudin (2019) on COMPAS

Now, take COMPAS – and its choice of **models**, in practice (Rudin 2019)

Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use
Interpretable Models Instead

Cynthia Rudin
Duke University

(i) Corporations can make profits from the intellectual property afforded to a black box.

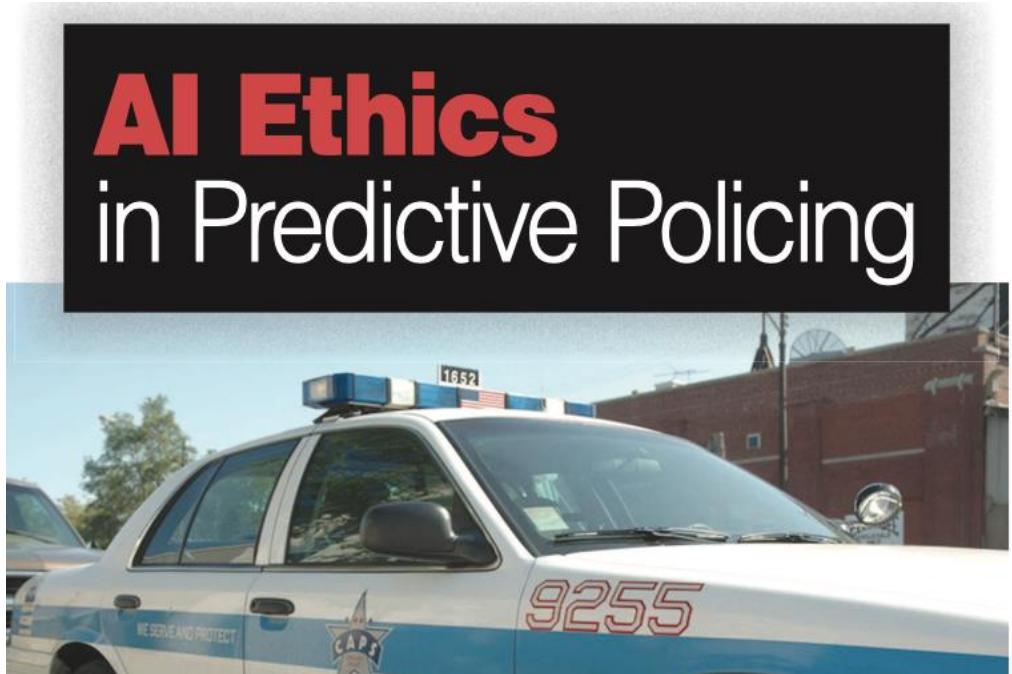
Companies that charge for individual predictions could find their profits obliterated if an interpretable model were used instead.

Consider the COMPAS proprietary recidivism risk prediction tool discussed above that is in widespread use in the U.S. Justice System for predicting the probability that someone will be arrested after their release [29].

The COMPAS model is equally accurate for recidivism prediction as the very simple three rule interpretable machine learning model involving only age and number of past crimes shown in Figure 3 below. However, there is no clear business model that would suggest profiting from the simple transparent model. The simple model in Figure 3 was created from an algorithm called Certifiably Optimal Rule Lists (CORELS) that looks for if-then

Asaro (2019) on PredPol (ref: COMPAS)

In the same vein we consider COMPAS;
cf PredPol in Asaro (2019)



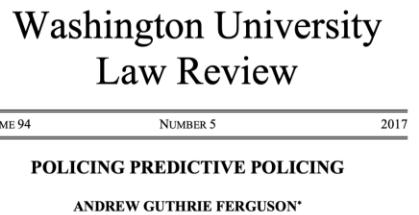
Transparency over the algorithms, data, and practices of implementation are also necessary. While the Chicago Police Department sought to avoid embarrassment from releasing the details of the SSL, it would be impossible for independent outside researchers to evaluate its impacts — positive and negative — without access to the data and algorithms. It should not take a prolonged lawsuit from a newspaper for government agencies to share public data. Of course, as more and more commercial systems, like PredPol,¹⁹ make the algorithms and even the data proprietary, they will fall under intellectual property protections. This means private companies will be processing the data, and will not be required to reveal their algorithms, or subject them to independent outside scrutiny. In some cases, private



Ferguson (2017) on PredPol (ref: COMPAS)

In the same vein we consider COMPAS; cf. PredPol - Ferguson (2017) provides a legal perspective.

Issues at all stages: crime stats; data dossiers; personal/cultural bias; data entry/analysis; tech complexity; financial/IP interests; auditing; metrics...



1. Transparency: Vulnerabilities

As currently implemented, a lack of transparency exists at all levels of predictive policing. Even something as simple as crime statistics, which in many cases are publicly available, remain rife with concerns about accuracy and completeness.³¹⁶ Adding personal data dossiers to these crime statistics creates new problems, as the sheer volume of information complicates transparent assessment of the sources underlying the predictions.³¹⁷ How do you fix an error in the data if you cannot see that such an error exists? How do you even know who has the responsibility to input information into these big aggregated databases?³¹⁸ In addition, unintended personal or cultural biases can infect the data, the scoring systems, the source codes, and thus the resulting predictive outcome.³¹⁹ Simply stated, without significant investment in exposing the data collection methods, weaknesses, and gaps, and without equal investment in understanding the challenges associated

with inputting and analyzing the data, the entire system runs the risk of being built on an unknown and unknowable database.³²⁰

The nature of algorithms further obscures the process, except perhaps to technical experts. Police officers and administrators receive the results, but due to the complexity of the chosen algorithm they can rarely understand the underlying math. Thus, predictive policing runs into the same problems as other automated predictive technologies: the technical complexity of the design makes it nearly impossible for outsiders to determine the accuracy, effectiveness, or fairness of the program.³²¹ True, police can see if the system works, but police cannot see how the system works. This lack of transparency is not simply the result of new technology, but also the influence of the proprietary nature of the software. The companies involved in these real-world tests are in a multimillion-dollar race to convince police departments to adopt their particular products. The companies have financial interests and proprietary secrets to protect, and every incentive to report positive outcomes.³²²

Effectiveness itself remains a contested issue. Early tests show a correlation between use of certain predictive policing techniques and decreased crime rates (for some crimes). But how do police districts determine metrics in the future? Crime may go up or down independent of the chosen computer program. Crime analysts may make a more or less accurate comparative judgment. Most importantly, how can outsiders audit the data? In similar police data collection experiments (DNA databases, “stop and frisk” reporting), the police have audited themselves with mixed results.³²³



Reflection.

AI systems in justice/policing causes serious effects on people's freedom and status under the law.

The legal perspective of AI ethics gives us another perspective on the need for transparency.

"Transparency is difficult, but it matters to a functioning predictive system that deals with individuals' lives and liberty"
(Ferguson, 2017).

How do we start fixing the issues?

E.g. for PredPol - auditing, public release of metrics, training
(Ferguson, 2017).





THE UNIVERSITY OF
MELBOURNE



Big Data Research & Social Media: *From Elections to Pandemics*



Reading: Walsh (2019)

Experiments in Social Media

Author: Toby Walsh¹

Abstract.

Social media platforms like Facebook and Twitter permit experiments to be performed at minimal cost on populations of a size that scientists might previously have dreamt about. For instance, one experiment on Facebook involved over 60 million subjects. Such large scale experiments introduce new challenges as even small effects when multiplied by a large population can have a significant impact.

Recent revelations about the use of social media to manipulate voting behaviour compound such concerns. It is believed that the psychometric data used by Cambridge Analytica to target US voters was collected by Dr Aleksandr Kogan from Cambridge University using a personality quiz on Facebook. There is a real risk that researchers wanting to collect data and run experiments on social media platforms in the future will face a public backlash that hinders such studies from being conducted. We suggest that stronger safe guards are put in place to help prevent this, and ensure the public retain confidence in scientists using social media for behavioural and other studies.



Reading: Walsh (2019)

Experiments in Social Media

Author: Toby Walsh¹

Abstract.

Social media platforms like Facebook and Twitter permit experiments to be performed at minimal cost on populations of a size that scientists might previously have dreamt about. For instance, one experiment on Facebook involved over 60 million subjects. Such large scale experiments introduce new challenges as even small effects when multiplied by a large population can have a significant impact.

Recent revelations about the use of social media to manipulate voting behaviour compound such concerns. It is believed that the psychometric data used by Cambridge Analytica to target US voters was collected by Dr Aleksandr Kogan from Cambridge University using a personality quiz on Facebook. There is a real risk that researchers wanting to collect data and run experiments on social media platforms in the future will face a public backlash that hinders such studies from being conducted. We suggest that stronger safe guards are put in place to help prevent this, and ensure the public retain confidence in scientists using social media for behavioural and other studies.

Not just Cambridge Analytica – which was for election targeting etc.

Walsh (2019) found that studies by academics “to improve voter participation” in fact “increased turnout by about 340,000 additional votes” (citing Bond et al, 2012).

Questions:

1. CA was bad, I’m sure you agree...
2. But for the 2nd experiment - isn’t this a good thing? Increasing voter participation = healthy democracy?

Reading: Walsh (2019)

The first recommendation is that we may need to take into account not just the impact on the individual under study but the broader impact any experiment might have on society. For a study on voting, this might be an electoral risk. For a study on fake news, it might be decreasing trust within society in real news. For a study on manipulating people's emotions, it might be the emotional wellbeing of the population studied.

Provocation/Thought
Experiment for #2 and #3:

A study on Twitter, say, might involve 100k – 1M of tweets/users (or more?)

How is this possibly achieved?

Reflections for transparency in social media and big data research, by Walsh (2019):

1. “The first recommendation is that we may need to take into account not just the impact on the individual under study but the broader impact any experiment might have on society...”
2. “The second recommendation is that ethics approval may be needed...”
3. “The third recommendation is that subjects of any experiment may need to be informed directly after the study about the results and their participation...”

Conclusion.

Big data research on social media invokes many concerns – privacy (can the user opt out of the ‘researchers gaze’); autonomy (does the research make the users do things they won’t otherwise?); wellbeing (does the research have the potential to change mood/health outcomes?)... **are these clear to the users?**

Concluding food for thought: is this ethical?

How Facebook and Google Track Public's Movement in Effort to Fight COVID-19

Location data provide rich resource for decision makers, scientists, and the public

By Emily Waltz



COVID-19 Community Mobility Report

Victoria 20 March 2021

Mobility changes

This data set is intended to help remediate the impact of COVID-19. It shouldn't be used for medical diagnostic, prognostic or treatment purposes. Nor is it intended to be used for guidance on personal travel plans.

The data shows how visits to places, such as corner shops and parks, are changing in each geographic region. Learn how you can use this report in your work by visiting [Community Mobility Reports Help](#).

Location accuracy and the understanding of categorised places varies from region to region, so we don't recommend using this data to compare changes between countries, or between regions with different characteristics (e.g. rural versus urban areas).

We'll leave a region out of the report if we don't have statistically significant levels of data. To learn how we calculate these trends and preserve privacy, read [About this data](#).

Sources: <https://spectrum.ieee.org/the-human-os/telecom/wireless/facebook-google-data-publics-movement-covid19>; Google



THE UNIVERSITY OF
MELBOURNE

Thank you

