

Renegotiating the Faustian Bargain for Data

Venkatesh Hariharan
June 2023

In a world where decisions are increasingly made inside algorithmic black boxes, how can people reclaim control of their data?

For people today, digital data has become a Faustian bargain. In German folklore, Dr. Faustus was someone who agreed to surrender his soul to the devil in exchange for worldly benefits. Similarly, people today surrender their personal data in exchange for services that are ostensibly free. The terms of this modern-day bargain, however, are even more opaque than the one of folklore. Whereas Dr. Faustus – disillusioned with life and the limited scope of human knowledge – made an intentional choice to enter into his deal, most people today share their data with little insight into and no control over how it is being used, collated, mined, and sold to the highest bidder.

Around the world, governments have proposed privacy laws to correct this situation, but getting these laws approved and setting up institutions like data protection authorities (DPAs) to enforce them can take years, if not decades. Even where privacy laws are being passed, these laws often provide significant exceptions to governments or the institutions necessary for oversight and enforcement of said laws are weak or nonexistent, thereby exacerbating the risks that data will be misused or abused.

Whether it is for digital services offered by the private sector or the public sector, the reality is that clicking the “I Agree” button on the incomprehensible end user license agreements (EULAs) usually leads to a permanent loss of data privacy. As Anja Kovacs of the Internet Democracy Project points out, “There is no other fundamental human right that we can sign away with the click of a button.”¹ With the rules of the digital playground heavily stacked against individuals, our choice is to either accept the situation or opt out completely and become a digital hermit.

Current privacy laws are of little help. In his paper, “Data Is What Data Does: Regulating Use, Harm, and Risk Instead of Sensitive Data,”² Daniel Solove argues that regulating data based on its type makes little sense in an age of big data and AI. Privacy laws in countries around the world have attempted to classify data into classes like personal data and sensitive personal data, with heightened protections for sensitive personal data. Solove convincingly argues that this approach is outdated in a world where big data and AI can be leveraged to make inferences about sensitive data.

Joanna Redden of the Data Justice Lab³ at Cardiff University explains how big data can harm us. In a *Scientific American* article titled, “The Harm That Data Do,”⁴ Redden writes that researchers studying the financial crash of 2008 found that banks had combined offline and online data to categorize and influence customers. As a result, the U.S. Department of Justice reached a \$175 million settlement with Wells Fargo over allegations that it had systematically pushed Black and Hispanic borrowers into more costly loans.

In our current industry structure, where institutions are data controllers and individuals are data subjects, the uses, harms, and risks of data are many. The Data Justice Lab highlights six major classes of big data harms, including targeting the vulnerable, misuse of personal information, discrimination, data breaches, political manipulation, social harms, and data and system errors. The lab maintains a Data Harm Record⁵ that is a running record of harms.

“We have entered an ‘age of datafication’ as businesses and governments around the world access new kinds of information, link up their data sets, and make greater use of algorithms and artificial intelligence to gain unprecedented insights and make faster and purportedly more efficient decisions,” writes Redden. “We do not yet know all the implications. The staggering amount of information available about each of us, combined with new computing power, does, however, mean that we become infinitely knowable—while having limited ability to interrogate and challenge how our data are being used.”

One of the difficulties of focusing on use, harm, and risk is that there are no frameworks for tracing

the chain of causality from use of the data to the harms caused. With most data processing happening in ways that are opaque to the user, tracing the link between data usage and harms to individuals remains a very difficult task. For example, users have little idea of how social media platforms, e-commerce systems, ride hailing apps, and others use their data. Worse, some websites and apps share data that users have not agreed to.

The State of Privacy 2022,⁶ a study by Arrka, an Indian privacy consulting firm, found that while 42% of Indian apps declare that they collect exact location data, the reality was that 76% of apps were collecting such sensitive data. The problem is compounded by the fact that data controllers suffer data breaches that compromise user data. For example, the author tried the Have I Been Pwned?⁷ website to see if his email and login had been compromised and found that 20 sites that he had registered at had leaked his data. No wonder that users feel they have very little control over their data. They might as well scatter feathers in the middle of a marketplace and try to gather them back.

By some estimates, around 3.5 quintillion bytes of data are generated every day as we browse the internet, online shop, message our friends and colleagues, watch online videos, and scroll through online publications. With the internet of things (IoT) connecting our toasters, refrigerators, and other devices to the internet, this amount of data is set to explode even further.

The digital footprints that we leave behind can be used to profile us in surprising ways. For instance, in a well-documented anecdote, a credit card company reduced the credit limit of a customer from \$10,800 to \$3,800, despite his high individual credit rating. The reason: This individual had used his card at the same store as shoppers who had poor repayment records.

Getting redress for such harms is tough. For more than a month, this customer tried to get in touch with the credit card company but did not receive an explanation for why his credit limit was cut. His case illustrates the kind of harm that can be caused by inferences made using big data. The Data Harm Record documents numerous ways in which big data profiling hurts individuals, including discrimination against qualified African-

American and Hispanic borrowers in mortgage lending, data brokers selling lists of financially vulnerable people, algorithmically based discriminatory pricing for SAT tests determined according to the location where people live, and minority neighborhoods paying 30% more for car insurance than white neighborhoods with the same risk levels.

Solove points out that, “To be effective, privacy law must focus on use, harm, and risk rather than on the nature of personal data. The implications of this point extend far beyond sensitive data provisions. In many elements of privacy laws, protections should be based on the use of personal data and proportionate to the harm and risk involved with those uses.”

But this is easier said than done. As Solove says, “Regulating based on use, harm, and risk is a difficult road, fraught with complexity, so it is no surprise it is often the road not taken.” At the regulatory level, policymakers will have to use various strategies like enforcing data fiduciary norms on data collectors, imposing fines and other deterrents like jail terms for violators, and ensuring algorithmic explainability (AE). However, lawmaking is a long, drawn-out process at the best of times, and these regulatory changes might take years to materialize.

Faced with permanent loss of control over data, many individuals have started taking defensive measures. This includes minimizing the sharing of personal photos on social media, requesting not be tagged on social media or disabling the ability to be tagged, turning off location data to minimize the risk of kidnapping, and selecting privacy friendly alternatives like the Mozilla Firefox browser and the DuckDuckGo search engine. While data sharing controls provided by digital services to their users is a small and welcome step in the right direction, the fact remains that users have to operate under the rules defined by these services. And even with privacy friendly services, users are constantly looking over their shoulders to see how their data trails are being followed and the kinds of inferences that are being made about them.

Still, people can’t spend their lives playing a defensive game over the data that belongs to them. As we become an increasingly digitized society, the task of managing one’s data should

not feel like a Sisyphean ordeal. Researchers like Kovacs argue that the status quo stems from an understanding of data as a resource to be mined and exploited. However, the Internet Democracy Project’s work has highlighted that such descriptions of data often do not match people’s experiences. Kovacs points out that victims of the nonconsensual sharing of sexual images generally do not describe the harms they experience in terms of a data protection or even privacy violation. Rather, they describe the harms as similar to those arising from sexual assault—a violation of bodily integrity.

“Taking people’s experiences as the starting point, thus, makes evident that in practice, the line between our physical and virtual bodies is increasingly becoming irrelevant—so much so, in fact, that maintaining the distinction is becoming harmful. In the digital age, bodies and data are closely intertwined. The nature and impact of data and data practices is embodied.”⁸

Using feminist theories and learnings around sexual consent, Kovacs argues that consent is not a simple Yes/No answer but the beginning of a process. “Consent has to be asked again and again for different situations and different acts. This really brings into picture the question of power relations.” Through her research, Kovacs points out that one of the most important aspects of privacy is that of boundary management.

Solove’s paper shines a light on where the focus of privacy law should be, while Redden and the Data Justice Lab have highlighted the real harms emerging from big data, and Kovacs explains the mental models that lead to such abuse. We are now entering a new era where legal methods alone are insufficient to protect privacy. Often, it is like locking the stables after the horses have bolted. We are entering an era of techno-legal regulations where technology and law have to work hand in hand. Fortunately, there are several emerging technologies that can help individuals renegotiate the Faustian bargain for their data.



Escaping the status-quo

In our current industry structure, where data is harvested through broad consent agreements and processed through big data, boundary management is a distant dream. Confronted with rampant encroachment on their privacy, individuals are looking for better ways to set boundaries and enforce them. A few emerging technologies could help, such as self-sovereign identity (SSI), verified credentials (VCs), account aggregators (AAs) and federated learning.

Self-sovereign identity: Identity was one of the missing pieces in the design of the internet. The work-around for this was usernames and passwords, which were followed by third-party login services provided by companies like Google, Facebook, Twitter, and LinkedIn. All of us have struggled to remember our usernames and passwords and suffered being locked out of important services like our bank accounts. While third-party login services offer the convenience of not having to remember usernames and passwords, they are akin to scattering breadcrumbs all over the internet. In return for convenience, these services track and profile your actions online.

In the physical world, trusted IDs like national ID cards, drivers licenses, and voting cards have primarily been issued by governments. The advantage of physical ID cards for individuals is that it is much harder to aggregate and profile them compared to digital IDs. Self-sovereign identity is an emerging technology that can be better than physical ID cards. SSI gives users complete control over how their identity and digital footprints are stored and handled. Using SSI, people can specify how much information they release to websites and apps. Since they control their information using their private keys, hacks to a website will not compromise their usernames and passwords.

One of the big advantages of SSI for people is that their mobile phone can serve as a digital wallet that stores identity and personal information. This enables them to keep personal data at their fingertips and present identity proofs when they are needed for verification. For example, if

someone is buying alcohol and needs to prove that they are 18, they can use SSI. This is more secure than handing over a physical identity card that reveals the user's name, date of birth, address, and other details.

In an article in Coinbase, Christopher Allen, a standards and identity practice specialist, argues that we need SSI now because governments and companies are sharing an unprecedented amount of information—cross-correlating information such as user viewing habits, purchases, where people are located during the day, where they sleep at night, and with whom they associate.⁹ The Bhutanese Government¹⁰ and the U.K. National Health Services are among the early adopters of SSI.

Verified credentials: SSI and verified credentials (VCs) are part of the same class of decentralized technologies. VCs are digital versions of credentials that people can present to parties that need them for verification.

While SSI deals with identity, VCs capture the trail of credentials that can be associated with that identity and give the identity owner greater control over their data trails. For example, a seller on an e-commerce platform might have accumulated a great brand reputation based on thousands of customer ratings and reviews across many years. However, control of that data remains with the platform and not with the individual seller. If the seller wants to migrate that data onto another e-commerce platform, it is hard to do. Decentralized e-commerce systems like the Open Network for Digital Commerce (ONDC)¹¹ and the Kochi Open Mobility Network (KOMN)¹² being pioneered in India will allow sellers to control their data trails and share them with a host of buyer and seller apps, instead of being captive to a centralized platform.

VCs and SSIs use a decentralized system called a trust triangle consisting of three parties: the issuer, the user, and the verifier. The issuer might be a university that issues a graduation certificate that the user presents to a potential employer who verifies it. These technologies benefit all three parties. For issuers, it can reduce incidence of frauds. Search for “fake university degrees” online and you get hundreds of sites offering to provide them. VCs enable issuers to issue certificates that can be easily verified online by potential

employers and others. Users benefit because they have complete control over how much data they share and with whom they share it. They can also revoke access to this data. Verifiers like employers benefit because they can verify VC-based certificates immediately instead of spending weeks and months waiting for the university to get back to them with a confirmation. It also helps them reduce their compliance burden and risks of data breaches as they need not store copies of the certificates on their IT systems.

VCS also help individuals share their resumes with potential employers, who can verify the candidate’s employment history quickly. This makes the process easier for everyone, and saves time and effort spent on verifying employment histories. VCs also make the process of applying for home loans, payday loans, and other financial products easier as they are machine verifiable and can be processed quickly by lending institutions. As we move away from closed loop ecosystems (CLEs) like centralized e-commerce platforms to open loop ecosystems (OLEs) like ONDC and KOMN, VCs will become increasingly important.

Account aggregator: The Data Empowerment and Protection Architecture (DEPA) is an ambitious attempt to rearchitect data flows from the current organization-centric model to an individual-centric model. The account aggregator

model introduced in India is one of the first implementations of DEPA and gives individuals more control over and insight into how their data is used. Though the AA framework is sector agnostic, its first application is in the world of finance.

Individuals can choose to download any AA app and link to their banks, provident funds, mutual funds, and insurance accounts. At present, 5.5 million consent requests have been fulfilled. The AA model consists of financial information users (FIUs), financial information providers (FIPs), the AAs, and the users themselves.

Once an individual has signed up, they can easily share their data with lending institutions, wealth managers, and other FIUs. This reduces the time and effort that individuals spend collecting paper documents from multiple sources like banks.

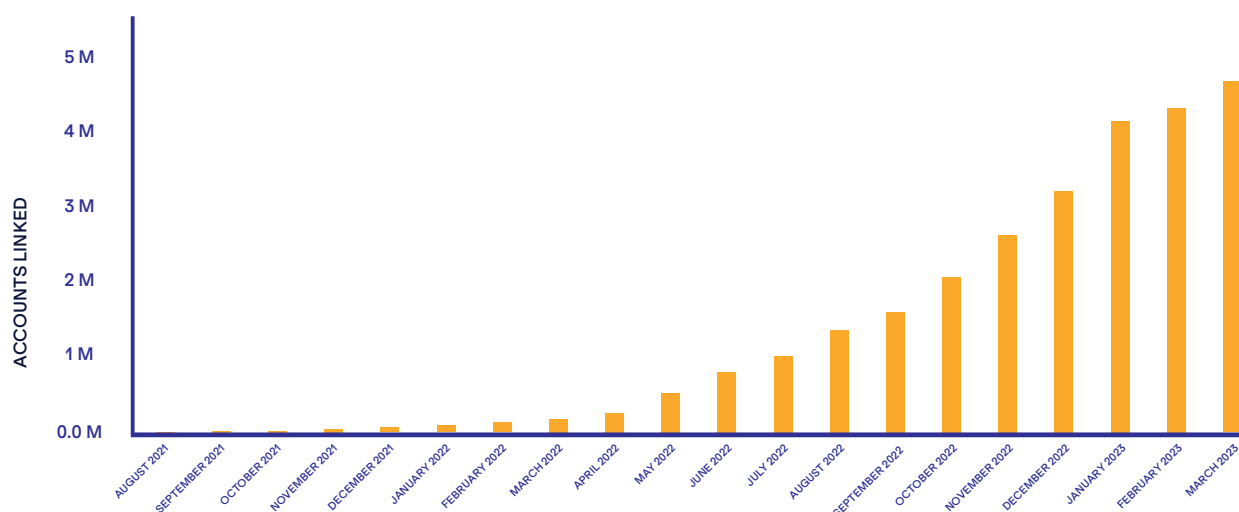
DEPA, and by extension AA, has been designed using the ORGANS framework, which is as follows:

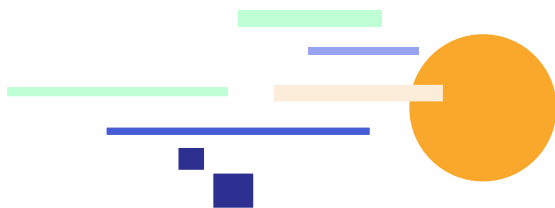
Open standards: The consent architecture must follow open standards and ensure all institutions use the same approach.

Revocable: The consent granted by the user can be revoked at any stage.

Source: www.sahamati.org.in

AA Ecosystem Performance Dashboard—FIP: Cumulative count of accounts linked by account-owners





Granular: Consent given has to be presented at a granular level, where the data is broken down in terms of its characteristics and how long it can be used, etc.

Auditable: All the events in the consent flow must be digitally signed and logged using the Ministry of Electronics and Information Technology's log artifact.

Notice: The user must be informed and given due notice when consent is created or revoked and when data has been requested, sent, or denied.

Security by design: The internal and external software to be used in DEPA must be designed from the ground up to be secure and provide end-to-end data security.

The ORGANS framework gives individuals an unprecedented level of control over their data. Apart from this, the AAs cannot monetize user data as the data flows between FIPs to FIUs are encrypted. Eventually, as more and more data streams become part of this ecosystem, the AA will give individuals greater control of their data sources that are scattered and difficult to control using current technology frameworks. A well-implemented consent network does the following:

- Gives individuals greater control over their digital data residing with multiple entities, like government departments, banks, mutual funds, hospitals, health care providers, and others, thereby enabling individuals to approve/reject data requests, revoke access to data, and share data at a granular level. For example, individuals can apply for a loan on an app, get information requests from multiple lenders, select the best lenders, and share data with them. This is more secure than screen scraping, where users need to share their accounts, customer data, and passwords with lending institutions. A consent network reduces the need for high levels of trust and the risk of security breaches that come with screen scraping.

- Shifts the data economy from an organization-centric architecture to an individual-centric one.
- Enables greater efficiencies in an economy by reducing friction in transactions. For example, a health care provider can use a consent network to access a patient's previous blood test reports stored across multiple pathology labs by sending a request to the patient. In the financial sector, a lender can quickly get back to a borrower by requesting digitally signed information that helps them assess the borrower's ability to repay. This benefits a vast majority of small and medium enterprises that have cash flows but do not have collateral to offer.

Federated learning: The norm today is that people send their data to centralized digital services that use their AI models to extract inferences about them. These inferences are then monetized. What if this model is flipped? Instead of people shipping their data to service providers, the data resides with them, and the data requesting organization sends their models to them. These models are not the black box algorithms run by centralized platforms but approved models that can be audited using algorithmic explainability principles. This ensures that if an individual has provided access to their data for the purpose of taking a car loan, the relevant model is sent to them, and their data is not used for any other purpose.

An emerging technology called federated learning enables models to be trained in a manner that preserves privacy. Using federated learning, hospitals could collaborate with each other to build and train models without sharing sensitive data with each other. This approach benefits everyone because AI models need more data than can be provided by an individual hospital. The improved models trained on distributed datasets across multiple hospitals benefit both patients, who get better diagnostics, and hospitals, which can improve the accuracy of treatments and patient turnaround times. Currently, there is still an element of centralization because the models are centralized. However, work is underway to decentralize the models, too. Another benefit of federated learning is that since data is distributed across multiple locations, it is less vulnerable than centralized servers to large-scale hacks that can lead to identity theft and fraud.

Conclusion

It is an old truism that givers have to set boundaries because takers have none. Faced with reports of massive data breaches that compromise their privacy and blatant misuse of their data like the Cambridge Analytica/Facebook case, and concerns over government overreach, people are increasingly taking matters into their own hands and opting for privacy-enhancing technologies. The inevitable conclusion is that individuals have to set their own boundaries in order to regain control of their data.

As technologies like SSI, VCs, AA, and federated learning become mainstream, they can help individuals set boundaries against data overreach by private-sector and government institutions. To be clear, such technologies cannot act as a substitute for well-drafted and implemented privacy laws. However, when such laws are in place, these technologies can act as complements to privacy laws.

Of course, laws and technology themselves are insufficient. The very structure of industry needs to be changed. We need to move away from an industry structure where institutions are data

controllers and individuals are data subjects, to one where individuals are data principals and data controllers, and institutions (whether private or government) are merely data fiduciaries that are held accountable for using data in the best interest of data principals. We need to move away from a system where individuals signing up for a digital service suffer a permanent and irrevocable loss of control to one where individuals share their data for a clear purpose and a limited period of time, and have the power to revoke access. Without such an industry structure, individuals will always be at the bottom of the data food chain. Aligning industry, governments, and technology to protect individual privacy might take many years, but the current status quo is simply unacceptable. The alternative visions that place individuals at the apex of the data food chain require greater awareness among individuals; a new class of technologies; rearchitecting data flows; and changes in privacy laws that focus on use, harms, and risk. If there is anything that history has taught us, it is that industry and governments will give up control reluctantly. Therefore, individuals have to stand firm in their desire to renegotiate the Faustian bargain for data.

Endnotes

- 1 Interview of Anja Kovacs on Swaddle, “Can the feminist idea of consent change how we think about privacy online?” <https://www.instagram.com/tv/CR7935sJ0Tw/>
- 2 “Data Is What Data Does: Regulating Use, Harm, and Risk Instead of Sensitive Data,” Daniel J. Solove, https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4322198
- 3 Data Justice Lab website, <https://datajusticelab.org/>
- 4 “The Harm That Data Do,” <https://www.scientificamerican.com/article/the-harm-that-data-do/>
- 5 Data Harm Record, <https://datajusticelab.org/data-harm-record/>
- 6 The State of Privacy 2022, <https://www.arrka.com/2023/01/26/the-arrka-study-2022/>
- 7 Have I been Pwned? <https://haveibeenpwned.com/>
- 8 “Towards an embodied approach to data,” <https://internetdemocracy.in/bodies-and-data>
- 9 “The path to self-sovereign identity,” <https://www.coindesk.com/markets/2016/04/27/the-path-to-self-sovereign-identity/>
- 10 “Bhutanese Crown Prince is country’s first citizen with digital identity,” <https://royalcentral.co.uk/asia/bhutanese-crown-prince-is-countrys-first-citizen-with-digital-identity-186659/>
- 11 Open Network for Digital Commerce, www.ondc.org
- 12 Kochi Open Mobility Network, <https://openkochi.net/>