



April 8, 2020

The Honorable Roger Wicker, Chairman
The Honorable Maria Cantwell, Ranking Member
U.S. Senate Committee on Commerce, Science, and Transportation
512 Dirksen Senate Office Building
Washington, DC 20510

RE: Using Privacy-Enhancing Technologies to ‘Enlist Big Data in the Fight Against Coronavirus’

Dear Chairman Wicker and Ranking Member Cantwell,

We write in support of your initiative to examine how public health responses to COVID-19 can respect fundamental privacy rights and data ethics in the upcoming hearing on ‘Enlisting Big Data in the Fight Against Coronavirus.’¹ This national emergency calls for privacy oversight that is vigilant in curbing the indefinite collection, retention, and repurposing of data; yet, adaptable to a health crisis that demands answers from data. Inpher submits the following comments to advise on the critical role of advanced, cryptographic privacy-enhancing technologies (PETs) in guiding accountable data-driven measures against COVID-19.

In an op-ed published on April 6 by Morning Consult (attached),² we illustrated why “lives over privacy” is a dangerous soundbite with an untenable premise. With the widespread growth of modern technologies such as Homomorphic Encryption and Secure Multi-Party Computation which can reconcile data utility with data privacy—time has come to reject the outdated assumption that deriving value from data requires a privacy tradeoff.

We strongly urge the Committee to consider, alongside legal and regulatory requirements, the need for built-in, technical safeguards for privacy that can help optimize beneficial uses of data in epidemiology whilst protecting individual and societal privacy against preventable risks.

Doing so would bring forward the U.S. privacy regime into the interdisciplinary framework exemplified by notable intergovernmental institutions. The European Data Protection Board (EDPB) which oversees the General Data Protection Regulation (GDPR),³ United Nations (UN),⁴ Organization for Economic Co-

¹ *Enlisting Big Data in the Fight Against Coronavirus*, 116th Cong. (2020), Senate Committee on Commerce, Science, and Transportation (Apr. 9, 2020), <https://www.commerce.senate.gov/2020/4/enlisting-big-data-in-the-fight-against-coronavirus>

² Sunny Seon Kang, *Privacy Does Not Pause in Pandemics*, Morning Consult (Apr. 6, 2020), <https://morningconsult.com/opinions/privacy-does-not-pause-in-pandemics/>

³ European Data Protection Board, *Guidelines 4/2019 on Article 25 Data Protection by Design and by Default* (Nov. 13, 2019), https://edpb.europa.eu/sites/edpb/files/consultation/edpb_guidelines_201904_dataprotection_by_design_and_by_default.pdf

⁴ United Nations, *UN Handbook on Privacy-Preserving Computation Techniques*, <http://publications.officialstatistics.org/handbooks/privacy-preserving-techniques-handbook/UN%20Handbook%20for%20Privacy-Preserving%20Techniques.pdf>



operation and Development (OECD),⁵ European Union Agency for Cybersecurity (ENISA),⁶ World Economic Forum (WEF),⁷ and World Health Organization (WHO)⁸ have uniformly endorsed the implementation of PETs to reduce systemic risks caused by an over-reliance on paper protections for privacy.

Inpher Background

We are a US-based cryptography and machine-learning company with the conviction that encryption and privacy are foundational to the future of computing and commerce. Inpher applies years of academic research⁹ on Fully Homomorphic Encryption (FHE) and Secure Multi-Party Computation (MPC) into the production of privacy-preserving analytics in healthcare, financial services, AI development, and more.¹⁰ Our technology, Secret Computing, builds on proprietary advances in privacy-enhancing cryptography to enable computing on encrypted and distributed datasets without revealing the underlying information across data sources.

Inpher's legal and policy department facilitates public education on privacy-preserving technologies, and promotes data protection by design and algorithmic accountability. We have testified before the U.S. House Committee on Financial Services on the utility of PETs in eliminating centralized security risks in third-party cloud services and customer data repositories.¹¹ We also consistently advocate for federal consumer protection agencies to impose systemic and technical privacy baselines that standardize PETs.¹²

Overview of Cryptographic Privacy Safeguards

Cryptographic technologies can provide a solution to traditional tradeoffs in privacy and analytical precision (by contrast to differential privacy), and allow secure collaboration across data silos for greater coordination and scalability.

⁵ OECD, *Revised Guidelines on the Protection of Privacy and Transborder Flows of Personal Data*: "The Joint Proposal also incorporates various recent data protection measures, including information management strategies, employee training, and appointment of individuals who are responsible for an organization's data protection practices, codes of practice, audits, privacy enhancing technologies, and privacy impact assessments."
https://www.oecd.org/sti/ieconomy/oecd_privacy_framework.pdf.

⁶ European Union Agency for Cybersecurity (ENISA), *Privacy Enhancing Technologies 'Time to Adopt PETs'*,
<https://www.enisa.europa.eu/topics/data-protection/privacy-enhancing-technologies>

⁷ World Economic Forum, *The Next Generation of Data-Sharing in Financial Services: Using Privacy Enhancing Techniques to Unlock New Value* (Sep. 12, 2019), <https://www.weforum.org/whitepapers/the-next-generation-of-data-sharing-in-financial-services-using-privacy-enhancing-techniques-to-unlock-new-value>

⁸ World Health Organization, *Global Health Ethics: Big Data and Artificial Intelligence*,
<https://www.who.int/ethics/topics/big-data-artificial-intelligence/en/>

⁹ Inpher, *High-Precision Privacy-Preserving Real-Valued Function Evaluation*, Financial Cryptography and Data Security (Dec. 2018) 183-202. https://www.researchgate.net/publication/335482512_High-Precision_Privacy-Preserving_Real-Valued_Function_Evaluation

¹⁰ Inpher, *Case Studies*, <https://www.inpher.io/case-studies-1#case-studies>

¹¹ Inpher, *Inpher CEO Dr. Jordan Brandt testifies before the U.S. House Financial Services Committee on "AI and the Evolution of Cloud Computing"* (Oct. 22, 2019), <https://www.inpher.io/news/brandt-testimony-artificial-intelligence-and-cloud-computing>; Testimony available here: <https://financialservices.house.gov/uploadedfiles/hhrg-116-ba00-wstate-brandtj-20191018.pdf>

¹² Inpher, *Privacy Advocacy and Public Policy*, <https://www.inpher.io/privacy-advocacy-and-public-policy>

Advances in MPC and FHE can be applied to (1) transform personal data into random auxiliary numbers that independently reveal no personally identifying information, and (2) ensure that the analyst conducting the computation only sees the output of the function without access to the private inputs contributed by multiple parties across organizations, industries, or jurisdictions.

This capacity is critical for a timely and accurate public health response. We examine below the need for cryptographic privacy safeguards in current proposals to enlist big data to combat COVID-19: (1) the potential collaboration of tech companies and the U.S. government to anonymize and aggregate location data for contact-tracing, and (2) the development of predictive AI systems for diagnosis and disease monitoring.

Risks in Current Data-Driven Proposals to Tackle COVID-19

Tech companies such as Facebook and Google are reportedly considering how anonymized and aggregated location data collected from their users can assist the U.S. government in tracing and suppressing COVID-19 infections.¹³

Indeed, aggregate data can effectively analyze movement and behavioral patterns which may unlock critical and timely information about COVID-19 transmissions. But when private, social media data is enlisted to inform public health policies enacted by the government, we must question if anonymization and aggregation alone will suffice.

Despite assurances from the spokespeople of Facebook and Google that no data will be directly shared with the government,¹⁴ there has been no public engagement¹⁵ on the following:

- (1) **Technical standards** that will guide the anonymization, aggregation, and analysis of user data, and the;
- (2) Risks of outsourcing public health analyses to companies that already enjoy **consolidated power from the centralization of valuable data.**

The assessment of risks should be both legal and technical. We must acknowledge the regulatory blind spots in the collection and retention of user data, and the vulnerabilities of mainstream anonymization techniques which can subject individuals to unforeseen risks such as:

¹³ Tony Romm, Elizabeth Dwoskin & Craig Timberg, *U.S. government, tech industry discussing ways to use smartphone location data to combat coronavirus*, Washington Post (Mar. 17, 2020), <https://www.washingtonpost.com/technology/2020/03/17/white-house-location-data-coronavirus/>

¹⁴ Alfred Ng, *Zuckerberg: Facebook isn't giving governments data to track coronavirus spread*, CNET (Mar. 18, 2020), <https://www.cnet.com/news/zuckerberg-says-facebook-is-not-giving-governments-data-to-track-coronavirus-spread/>

¹⁵ Casey Newton, *The US government should disclose how it's using location data to fight the coronavirus*, The Verge (Mar. 31, 2020), <https://www.theverge.com/2020/3/31/21199654/location-data-coronavirus-us-response-covid-19-apple-google>

- **Re-identification risks that undermine current anonymization techniques.** In July 2019, research published by Imperial College London demonstrated high success rates in re-identifying individuals in “anonymous” and incomplete datasets using machine-learning models. Using a generative model, they found that 99.98% of Americans would be correctly re-identified in any dataset using 15 general demographic attributes. Even if all other personal identifiers besides broad demographic data (gender, date of birth, and zip code) were removed, individuals were re-identified 54% of the time.¹⁶ The paper concluded, “our results suggest that even heavily sampled anonymized datasets are unlikely to satisfy the modern standards for anonymization set forth by GDPR and seriously challenge the technical and legal adequacy of the de-identification release-and-forget model.”¹⁷
- **Limited datasets create bias and perpetuate socio-economic divides.** Social media platforms hosted by Facebook and Google certainly do have a wide userbase. However, it is not a foregone conclusion that this is necessarily the most *representative* dataset available. We cannot safely rely on mobile user data to portray an accurate cross-section of the populations affected by COVID-19.¹⁸ Location data collected through the use of these platforms may over-represent privileged majorities and under-represent marginalized minorities who are at the greatest risk of infection.

In ‘The Pandemic’s Missing Data,’¹⁹ The New York Times reported that health authorities are lacking access to demographic data that is central to understanding and rectifying injustice with equitable testing and treatment. This gap applies to both COVID-19 tracing efforts, and the development of AI for predictive diagnosis and disease monitoring. In both, relying on a singular source of aggregate data leaves out insights on large sects of the population without access to Wi-Fi, mobile phones, social media accounts, and healthcare—skewing the accuracy of predictions on the rate of transmission. Conducting analysis on datasets that remove sensitive demographic data points such as gender, age, ethnicity, and sexual orientation, may actually hinder informed public health responses that can specifically respond to the disproportionate impact of COVID-19.²⁰

¹⁶ Caroline Brogan, Anonymising personal data ‘not enough to protect privacy’, shows new study, Imperial College London News (Jul. 23, 2019), <https://www.imperial.ac.uk/news/192112/anonymising-personal-data-enough-protect-privacy/>

¹⁷ Rocher, L., Hendrickx, J.M. & de Montjoye, Y. Estimating the success of re-identifications in incomplete datasets using generative models. *Nat Commun* 10, 3069 (2019). <https://doi.org/10.1038/s41467-019-10933-3>, See also, Charlotte Jee, *You’re very easy to track down, even when your data has been anonymized*, MIT Technology Review (Jul. 23, 2019), <https://www.technologyreview.com/2019/07/23/134090/youre-very-easy-to-track-down-even-when-your-data-has-been-anonymized/>

¹⁸ Jacob Hoffman-Andrews & Andrew Crocker, *How to Protect Privacy When Aggregating Location Data to Fight COVID-19*, Electronic Frontier Foundation (Apr. 6, 2020), <https://www.eff.org/deeplinks/2020/04/how-protect-privacy-when-aggregating-location-data-fight-covid-19> (“Smartphone ownership remains a proxy for relative wealth, even in regions like the United States where 80% of adults have a smartphone”), Pew Research Center, *Mobile Fact Sheet* (Jun. 12, 2019) <https://www.pewresearch.org/internet/fact-sheet/mobile/>

¹⁹ Aletha Maybank, *The Pandemic’s Missing Data*, New York Times (Apr. 7, 2020), <https://www.nytimes.com/2020/04/07/opinion/coronavirus-blacks.html>

²⁰ Academy of Medical Royal Colleges, *Artificial Intelligence in Healthcare* (Jan. 2019), https://www.aomrc.org.uk/wp-content/uploads/2019/01/Artificial_intelligence_in_healthcare_0119.pdf (“The UK Government and its health and social care systems have a legal duty to maintain the privacy and confidentiality of its citizens...However, the development of AI and machine learning algorithms relies on the use of large datasets.”)

- **Worsening the information asymmetry for consumers.** Consent to activate location services within a social media app should not extend to accepting the risks of re-identification and targeting in a national pandemic registry. Facebook²¹ and Google,²² among other Big Tech companies, have a history of violating the Federal Trade Commission’s data privacy and security regulations against unfair and deceptive data practices. Empirical research has shown that these companies systematically employ “dark patterns” designed to deceive users into relinquishing more data than is necessary, and to discourage them from exercising privacy controls on the platform.²³ This means that once the data is collected, tech companies are insulated from public accountability by their proprietary access to the data.
- **Single point-of-failure inherent in centralized data systems.** Entrusting a few tech incumbents to collect, safeguard, and compute big data for COVID-19 tracing creates a data oligarchy. Beyond the obvious anti-competitive effects of this infrastructure, this creates a central entity that is highly susceptible to a single point-of-failure to data breaches and information abuse by malicious or simply negligent actors.

The Role of Privacy-Enhancing Cryptographic Techniques in Tracing and AI

All these risks have common remedies. Accurate data-driven responses to COVID-19 require more data points.²⁴ Equitable public health responses require more data sources to inform policymakers which vulnerable groups are left out of current efforts. Long-term privacy and security require the decentralization of data. Cryptographic privacy safeguards have the potential to address them all.

Collaborative information-sharing on advanced privacy-preserving technologies such as MPC and FHE will be critical for systemic accountability and data protection in a time of crisis, because:

- Distributed and encrypted computing can make sensitive data points accessible to multiple parties by transforming plaintext data into random auxiliary numbers, which independently reveal no identifying information.
- After each computing operation, these random auxiliary numbers (also called “triplets”) are deleted, thereby minimizing data and imposing a strict purpose limitation.

²¹ U.S. Federal Trade Commission, *In re Facebook, Inc.* (last updated Jul. 24, 2019), <https://www.ftc.gov/enforcement/cases-proceedings/092-3184/facebook-inc>

²² U.S. Federal Trade Commission, *Google and YouTube Will Pay Record \$170 Million for Alleged Violations of Children’s Privacy Law* (Sep. 4, 2019), <https://www.ftc.gov/news-events/press-releases/2019/09/google-youtube-will-pay-record-170-million-alleged-violations>; *In re Google, Inc.* <https://www.ftc.gov/enforcement/cases-proceedings/102-3136/google-inc-matter>

²³ Norwegian Consumer Council, *Deceived by Design* (Jun. 27, 2018), <https://fil.forbrukerradet.no/wp-content/uploads/2018/06/2018-06-27-deceived-by-design-final.pdf>

²⁴ Joseph Bullock, Alexandra Luccioni, et al, *MAPPING THE LANDSCAPE OF ARTIFICIAL INTELLIGENCE APPLICATIONS AGAINST COVID-19*, Cornell Computers and Society (Mar. 25, 2020), <https://arxiv.org/pdf/2003.11336.pdf>; Felicia Vacarelu, *Mapping the landscape of artificial intelligence applications against COVID-19*, United Nations Global Pulse (Mar. 26, 2020), <https://www.unglobalpulse.org/2020/03/mapping-the-landscape-of-artificial-intelligence-applications-against-covid-19/>

- Unlike classic anonymization techniques, advanced cryptography can unlock the analysis of sensitive data points with multiple data sources—without disclosing or transferring them, or having to factor them out with de-identification methods.

This de-centralized yet collaborative capacity brings remarkable public value to COVID-19 efforts:

- **AI is interdisciplinary in nature and needs access to heterogeneous datasets.** A joint paper by the UN Global Pulse and WHO affirms the need for privacy safeguards that can facilitate data collaboration in healthcare AI:²⁵

First, we believe that scalable approaches to data sharing using open repositories will drastically accelerate the development of new models and unlock data for the public interest. Image-based medical diagnosis, in particular, is a domain in which training data is currently scarce but the value of AI models may be high.

In order to facilitate the sharing of such data, clinical protocols and data sharing mechanisms will need to be designed and data governance frameworks will need to be put in place. It is important to reinforce that research with medical data must be subject to strong regulatory requirements and privacy protecting mechanisms. Overall, any AI application developed should undergo an assessment to ensure that complies with ethical principles and above all respects human rights.

- **Cryptographic privacy safeguards can enhance access to high-quality and traceable data** in AI models, and also maintain compliance with applicable privacy laws and security standards in the U.S. and internationally. The distributed and coordinated computing capacity of MPC also protects the dataset from a single point of security failure.
- **Cryptographic PETs are critical to institutional capacity-building for healthcare AI** in both private and public sectors. Multilateral and multimodal data sources in AI systems can ensure that the model is scalable to complexity, interoperable, representative, and ethically deployable.

Privacy and security experts around the world agree that we must move away from “omniscient central servers”—to quantifiable cryptographic methods that facilitate collaboration without the transfer of personal information.

- CCC (European Information Security Association): “A dependence of the users' privacy on the trustworthiness and competence of the operator of central infrastructure is technically not necessary. Concepts based on this ‘trust’ are therefore to be rejected.”²⁶

²⁵ *Id.*, at 10.

²⁶ Chaos Computer Club (CCC), *10 requirements for the evaluation of "Contact Tracing" apps* (Apr. 6, 2020), <https://www.ccc.de/en/updates/2020/contact-tracing-requirements>

- DP-3T, a group of academics from the Swiss Federal Institute of Technology (EPFL), KU Leuven, University College London, Oxford University, and others has published a call to action for ‘Decentralized Privacy-Preserving Proximity Tracing’:²⁷

Designs with centralized components, where a single actor, such as a server or a state, can learn a great deal about individuals and communities, need specific attention because if they are attacked, compromised or repurposed, they can create greater harm.

In order to address these issues, we instead realize the same task using a decentralized design that does not require the centralized collection and processing of information on users. Such a design builds on strong, mathematically provable support for privacy and data protection goals, minimizes the data required to what is necessary for the tasks envisaged, and prevents function creep, for example for law enforcement or intelligence purposes, by strictly limiting how the system can be repurposed with cryptographic methods.

Given the direct impact of private COVID-19 tracing on government policies, the public should not be left in the dark about the calculus of the necessity, proportionality, and lawfulness of these measures.

The U.S. does not have a federal privacy law to establish a strong foundation of accountability for proposed data-driven measures—which may have long-term consequences for privacy and civil liberties if compromised through re-identification, or abused with unlawful applications after the pandemic ends. Inpher strongly believes that in the absence of sufficient legal safeguards against data misuse, technical safeguards should be employed to automate data minimization, limited retention, and purpose limitation.

Conclusion

Our comments underscore the need for cryptographic privacy safeguards in the data ecosystem underpinning COVID-19 responses. This crisis is an impetus for privacy-enhancing technologies to guide evidence-based policymaking with respect to the collective, public value of privacy. We hope that our arguments and the materials we have cited will help the Committee in considering technical safeguards for privacy that can reduce barriers to innovative public health technologies that are efficient, accurate, and equitable.

Thank you for the opportunity to comment, and we kindly ask that this letter and the attachments be entered into the hearing record.

²⁷ *Decentralized Privacy-Preserving Proximity Tracing: Overview of Data Protection and Security* (Apr. 3, 2020) available at, <https://github.com/DP-3T/documents/blob/master/DP3T%20-%20Data%20Protection%20and%20Security.pdf>



Sincerely,

A handwritten signature in black ink, appearing to read "Sunny Seon Kang".

Sunny Seon Kang

Senior Privacy Counsel, Head of Policy

Inpher, Inc.

sunny@inpher.io

Attachments:

(1) World Economic Forum, *The Next Generation of Data-Sharing in Financial Services: Using Privacy Enhancing Techniques to Unlock New Value* (Sep. 12, 2019). Pages 13 – 19 discuss Homomorphic Encryption, Zero-Knowledge Proofs, and Secure Multi-Party Computation.

(2) Sunny Seon Kang, *Privacy Does Not Pause in Pandemics*, Morning Consult (Apr. 6, 2020)

White Paper

The Next Generation of Data-Sharing in Financial Services: Using Privacy Enhancing Techniques to Unlock New Value

Prepared in collaboration with Deloitte

September 2019



World Economic Forum
91-93 route de la Capite
CH-1223 Cologny/Geneva
Switzerland
Tel.: +41 (0)22 869 1212
Fax: +41 (0)22 786 2744
Email: contact@weforum.org
www.weforum.org

© 2019 World Economic Forum. All rights reserved. No part of this publication may be reproduced or transmitted in any form or by any means, including photocopying and recording, or by any information storage and retrieval system.

This white paper has been published by the World Economic Forum as a contribution to a project, insight area or interaction. The findings, interpretations and conclusions expressed herein are a result of a collaborative process facilitated and endorsed by the World Economic Forum, but whose results do not necessarily represent the views of the World Economic Forum, nor the entirety of its Members, Partners or other stakeholders.

Contents

Foreword	4
Introduction	5
Chapter 1: Privacy in the financial sector	6
The benefits of data-sharing	6
The potential drawbacks of data-sharing	6
Changing the dynamics of data-sharing	6
Chapter 2: Privacy enhancing techniques	8
Technique #1: Differential privacy	9
Technique #2: Federated analysis	11
Technique #3: Homomorphic encryption	13
Technique #4: Zero-knowledge proofs	15
Technique #5: Secure multiparty computation	17
Chapter 3: Applications in financial services	20
Unlocking new value for financial institutions	20
Unlocking new value for customers	22
Unlocking new value for regulators	24
Closing comments	25
Appendix	27
Benefits and limitations of techniques	27
Further reading	30
Acknowledgements	31
Endnotes	33

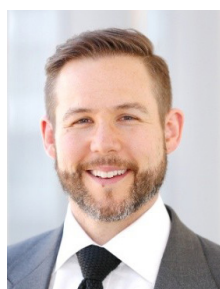
Foreword



Matthew Blake,
Head of Future
of Financial
and Monetary
Systems, World
Economic Forum



Jesse McWaters,
Financial
Innovation Lead,
World Economic
Forum



Rob Galaski,
Global Leader,
Banking & Capital
Markets, Deloitte
Consulting

The centrality of data to the transformations of the Fourth Industrial Revolution is today so self-evident as to have become a cliché, and whether you believe data is the new oil, the new gold or even [the new bacon](#), there is no doubt that its growing importance is shifting the priorities of the private sector. However, while many column inches have been dedicated to the competitive scramble to accumulate vast troves of data, less attention has been paid to the growing appetite of firms to unlock the power of data-sharing between institutions. Within the financial system specifically, we have seen a significant increase in the appetite for such collaborations across use cases ranging from improving fraud detection to enabling new forms of personal financial advice.

Of course, sharing data is not without risks. The potential value of collaboration must be weighed against its implications on customer privacy, data security and control of competitively sensitive data. Historically, this balance between privacy and utility has created tensions and conflicting objectives in the financial services industry, where any value obtained through data-sharing often needed to be weighed against the potential increase in privacy risks. These tensions have seen many seemingly promising opportunities for data-sharing shelved long before they could be deployed.

However, an emerging set of technologies called “privacy enhancing techniques” have the potential to fundamentally redefine the dynamics of data-sharing by eliminating – or greatly reducing – the risks historically associated with collaboration. As these technologies mature, they will demand a re-examination of a host of mothballed data-sharing projects and the exploration of previously unimaginable opportunities.

Privacy enhancing techniques have the potential to unlock enormous value for the financial sector – but they will do so only if senior executives and regulators have an awareness and working understanding of these mathematically and computationally complex techniques. The purpose of this paper is to provide an abstract and easy-to-grasp understanding of some of the most promising techniques emerging today and an illustration of how they might be deployed in the financial system. In doing so, we hope to support the emergence of a more collaborative financial environment where shared data can lead to shared benefits for financial institutions, customers and the broader financial system.

Introduction

It is an age-old tale – three blind men stumble upon an elephant for the first time. One feels its leg and concludes that the elephant is a tree. One feels the trunk and thinks the elephant is a large snake. The last feels its tail and surmises the elephant is a broom.

In the financial services sector, institutions face a similar challenge of not being able to “see the whole elephant”; each institution holds a piece of the puzzle (i.e. data) when it comes to answering important questions such as “is this customer creditworthy?”, “are these traders colluding?”, or “is this transaction fraudulent?” However, with only their own data, financial institutions – like the three blind men – risk drawing the wrong conclusions. In the parable, sharing information is the key to unlocking the mystery of the elephant and building a complete picture of the pachyderm at hand. Unfortunately, this kind of data-sharing is not so easy for financial institutions. Unlike the blind men, they face many restrictions on how they store, manage and share data that, until recently, have made it impossible for them to build a comprehensive picture of their customers and operating environments.

The value of the whole of data is greater than its component parts, but capturing this value is fraught with complexity and conflicting goals. For example, by sharing data, financial institutions would be able to better identify patterns that suggest transaction fraud, leading to fewer false positives in the detection of financial crime. However, they are wary of disclosing valuable competitive intelligence on their customer base, and of creating tensions with privacy regulations. It is important to note that it is not only financial institutions that stand to benefit: By sharing data, customers would be able to benefit from more personalized, specific and nuanced advice. However, they are wary of their data being misused, abused and shared without their consent.

These examples highlight the tensions of sharing data; there is value to be derived from doing so, but it traditionally diminishes privacy (of the individuals whose data is being shared) and confidentiality (of the institutions supporting the data-sharing). Historically, great effort has been dedicated to navigating these conflicting objectives and operating the financial system in a way that institutions, customers, civil society and regulators are all amenable to. “Privacy enhancing techniques” allow institutions, customers and regulators to unlock the value in sharing financial data without compromising on the privacy and confidentiality of the “data owners” (i.e. customers) and “data stewards” (i.e. financial institutions). These techniques are not new, but significant developments in recent years have transformed them from research curiosities to production-ready techniques with the potential to alter the fundamental nature of data-sharing.

This document is intended for use by executives at financial institutions across subsectors (e.g. insurance, banking, investment management); it provides a high-level overview of how these privacy enhancing techniques work, and the value they can unlock within financial institutions. In this White Paper, we will:

Chapter 1: Take a closer look at the tensions surrounding privacy in the context of the financial sector

Chapter 2: Understand how several privacy enhancing techniques work

Chapter 3: Demonstrate how they could be used to enable new types of data-sharing

This report include three chapters:



Chapter 1:
Privacy in the
financial sector

p. 6



Chapter 2:
Privacy enhancing
techniques

p. 8



Chapter 3:
Applications in
financial services

p. 20

Chapter 1: Privacy in the financial sector

Competing objectives surrounding the use of data pull financial institutions in a variety of different directions when it comes to deciding how data is to be stored, managed and shared. These tensions have historically existed across three different domains: within institutions themselves, their regulators and their customers.



Below, we explore the conflicting objectives (the benefits of data-sharing and its drawbacks) for each of these three domains.

The benefits of data-sharing

Financial institutions can benefit from three forms of data-sharing:

- Inbound data-sharing (acquiring data from third parties)
- Outbound data-sharing (sharing owned data with third parties)
- Collaborative data-sharing (inbound and outbound sharing of similar forms of data)

Inbound data-sharing allows institutions to enrich their decision-making systems with additional information, leading to higher-quality outputs and more accurate operations. For example, trading firms can use third-party services such as Thomson Reuters MarketPsych Indices¹ to inform their buy/sell decisions with social media data, hypothetically leading to a more accurate understanding of market sentiment. Outbound data-sharing, on the other hand, allows institutions to draw on capabilities (and offer customer benefits) that they may not own internally. For example, Wealthsimple, a robo-adviser, allows its clients' portfolio information to be pulled into Mint.com through a secure connection,² so that customers can see their investment balances alongside their day-to-day spending and build a comprehensive understanding of their finances. Finally, collaborative data-sharing allows institutions to achieve a scale of data that they would not be able to reach on their own, unlocking a depth and breadth of insights that would otherwise not be possible. For example, six Nordic banks recently announced a collaboration to develop a shared know-your-customer (KYC) utility³ that will allow them to strengthen their financial-crime prevention systems.

For regulators, data-sharing presents an opportunity to return control and ownership of financial data back into the hands of customers, ultimately leading to increased competition and innovation. This is seen in the Open Banking Standard in the UK, PSD2 more broadly in the EU, the Consumer Data Right in Australia, and other forms of Open API regulations in Singapore, Hong Kong and Japan. Each of these regulations, in some form, requires institutions to make the data they hold on their customers (e.g. transaction data) available to accredited third parties as requested by the customer. This

allows new market participants to access the data and build new value propositions; ultimately, regulators believe this will lead to improved financial outcomes for citizens.⁴

For customers, sharing data allows them to receive specific benefits – whether in the form of higher-quality products or more efficient services. For example, Lenddo provides customers with a higher-quality (i.e. potentially more accurate) credit score by analysing their social media data, telecom data and transaction data.⁵ Customers are increasingly aware of the value of their personal data and seek to share it (whether by directly providing an institution with more information or authorizing an institution to share their data with a third party on their behalf) only when the benefits received in exchange are meaningful.⁶

The potential drawbacks of data-sharing

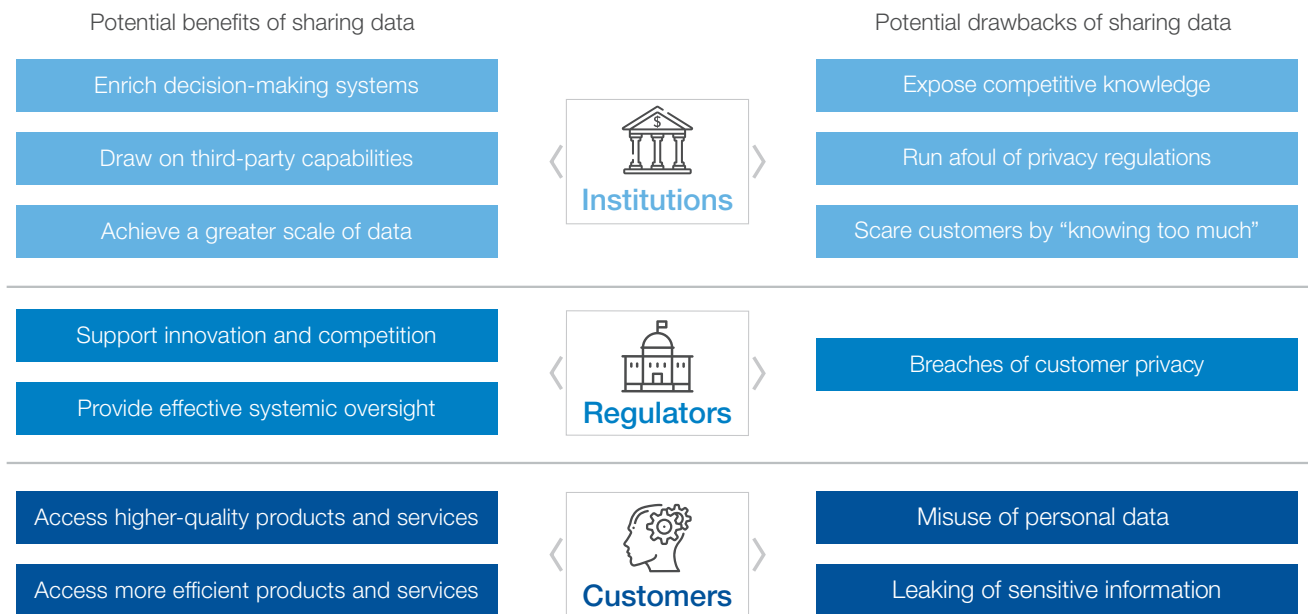
However, there are also several factors that inhibit the sharing of data in financial services. For financial institutions, any outbound data-sharing presents the risk of exposing competitive knowledge (e.g. the identities of customers and their characteristics) that could be misused by third parties. Furthermore, sharing data may run afoul of privacy regulations such as GDPR, or introduce complexities to the necessary processes (e.g. building out new mechanisms to ensure informed consent) that outweigh the potential benefits. And finally, with the increasing use of AI and other advanced analytical techniques, executives at large financial institutions have begun to worry about the “creep factor” – knowing too much about a customer and alarming them.

For regulators, protecting consumers' financial and non-financial confidentiality is a critical responsibility, and limiting the sharing of data has historically been the instrument to achieve this.⁷ For example, the Gramm-Leach-Bliley Act of 1999 in the United States requires financial institutions to communicate how their customers' sensitive data is being shared, allow them to opt out, and apply specific protections on what is shared.⁸ In recent years, regulatory authorities around the world have also introduced new and more stringent customer privacy requirements. For example, GDPR in the EU

requires institutions to, among other things, provide customers with easier access to the personal data about them held by the institution. Other regulations prevent firms from sharing personally identifying information (PII) across country borders to protect national customer privacy, potentially preventing multinational institutions from analysing their own internal data throughout their organization. Such requirements make certain types of data-sharing impossible, or so expensive, complex and time-consuming that the business case for doing so is weakened.

Finally, while customers seek additional benefits from sharing their data, they are also increasingly wary that their data could be misused by the firms that hold it: A survey conducted by Harris Poll shows that only 20% of US consumers “completely trust” the organizations they interact with to maintain the privacy of their data.⁹

This is no doubt exacerbated by several high-profile security and privacy breaches in 2018, including Cambridge Analytica,¹⁰ Capital One,¹¹ Google+,¹² Aadhaar¹³ and others. Customers fear that their data could be used to harm them (e.g. through identity theft) and more broadly that unintended parties can learn something about them that they wish to keep private (e.g. sensitive purchase history).¹⁴



Changing the dynamics of data-sharing

As illustrated, privacy tensions exist for every stakeholder in the financial services sector, and navigating these tensions has historically left significant value in data-sharing uncapturable. However, emerging privacy enhancing techniques are enabling institutions, customers and regulators to share data in a way that helps to achieve a balance between competing opportunities and obligations, allowing for data-sharing that is compliant with regulatory principles, protects the privacy of customers and safeguards the confidentiality of institutions’ business processes. These techniques have the potential to expand the range of feasible data-sharing opportunities in financial services, effectively allowing institutions to “see the whole elephant” and unlock new value for themselves, their customers, regulators and societies at large.

Chapter 2: Privacy enhancing techniques

Data acts as the fuel to the fire powering the Fourth Industrial Revolution, underpinning the growth of new technologies such as artificial intelligence and connected devices. To truly benefit from these new technologies, institutions need to be able to use the data available to them, both within their institutions and outside of them. Below, we outline five key techniques to managing data privacy that are enabling institutions to unlock new value. These five techniques¹⁵ are:



For each of these privacy enhancing techniques, we will: explore the potential benefits; demonstrate how they work with a hypothetical case; illustrate how they can be useful through historical cases of privacy failures; and assess the viability of the technique in financial services. We will then discuss how these techniques can be combined to enable new data-sharing collaborations in the industry.



Technique #1: Differential privacy

Overview:

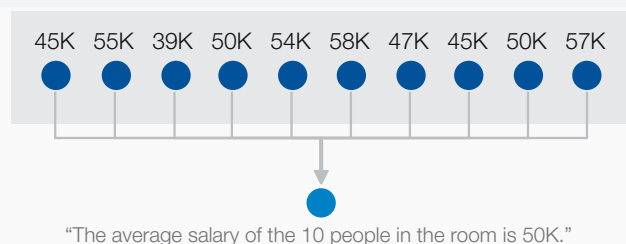
When an institution is seeking to share data with a third party, removing or anonymizing personally identifiable information is not always enough to protect the privacy of the individuals in the database. For example, the data could be correlated to other datasets to reidentify specific individuals in the database. One well-established way to address this is to add noise to the process (to the inputs, the calculations themselves or to the outputs), ensuring the privacy of individual “rows” of data while meaningful insights can still be derived from queries on the aggregate data. For example, census data is often anonymized with noise to protect the privacy of individual respondents; in the United States, differential privacy will be used for the 2020 Federal Census.¹⁶

In 2006, Cynthia Dwork et al.¹⁷ published a hallmark paper on “differential privacy”, providing a generally applicable mechanism to calculate the amount of noise that needs to be added to data to protect the privacy of every individual within the database.¹⁸ Since then, significant additional research has been advancing the efficiency and scalability of this approach, and it has been adapted into a variety of real-world applications. Differential privacy is used in large-scale production environments by companies such as Apple (e.g. to autocomplete web searches¹⁹) and has been embedded into a variety of popular analytics and machine-learning libraries such as PyTorch²⁰ and TensorFlow.²¹

Note: differential privacy is not a technique/mechanism itself, but a measurement of various techniques and methods of adding noise that limit the ability of an outsider attempting to deduce the inputs to an analysis from the results of the analysis.

How it works:

Consider a hypothetical case where a group of 10 individuals with the same job are seeking to share their salary information to understand if they are overpaid or underpaid, but they do not want to disclose their actual salary figures to any of the other individuals. In order to do so, they ask an independent and trusted third party to act as an intermediary, anonymizing their inputs while still providing useful insights on the aggregated data. The intermediary averages their data and informs them that the average salary of the 10 individuals in the room is 50K. This is useful information to the individuals as they can directionally establish whether they are overpaid or underpaid.

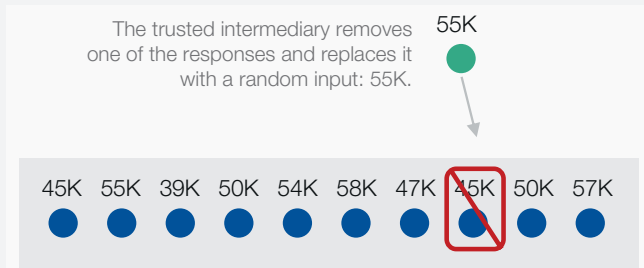


However, consider the case where one participant already has access to the salary data of eight others in the room, leaving only one individual’s information unknown.

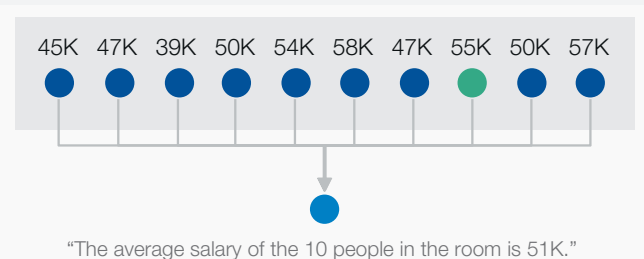


Once the average salary of the room is known, this individual can deduce that the exact salary of the 10th individual is 45K and can expose/use this private information.

To prevent this privacy breach, the intermediary could add noise to his/her calculation of the average. For example, the surveyor could remove one of the 10 participants’ responses, and replace it with a random number within the range of the maximum and minimum responses received (i.e. between 39K and 58K).



By then calculating the average salary as usual, the intermediary provides a slightly noisy response of 51K, and makes it impossible for any third party to reverse-engineer the inputs provided.



The individual with knowledge of the eight other salaries cannot deduce the exact salary of the final person in the room, since the process of adding noise creates two uncertainties:

- Any one of their eight known salaries may have been replaced by an unknown number, leading to a possible salary range of 36–74K for the unknown individual when an average of 51K is provided as the output. This range is so large that it provides no value.
- The unknown individual's salary itself was removed from the sample set, in which case not even a salary range could be reverse-engineered.

The individual seeking to breach the privacy of the respondents does not know which of the two above situations has occurred, and thus cannot reverse-engineer the salary information of the last individual in the room. Meanwhile, the others can still directionally ascertain whether they are overpaid or underpaid.

If the intermediary cannot be trusted to keep individuals' information private, they can also instead add noise to their individual inputs prior to sharing with the intermediary. For example, they can each add or remove up to a certain allowance (e.g. 2K) to the number they provide to the intermediary. The output will still be directionally correct and allow individuals to ascertain whether they are overpaid or underpaid, while protecting the privacy of their individual inputs.

Where it could have helped:

In the mid-1990s, a state government insurance body released anonymous health records to encourage public research in medical care. The data had been anonymized using several techniques, e.g. addresses had been removed and names had been replaced with randomized strings. However, researchers were able to compare and correlate this information with publicly available voter registration data to reidentify many individuals in the database,²² including state officials who had previously assured the public that patient privacy was protected. Rather than exposing the database directly, a differential privacy system could be implemented to take queries on the dataset and add noise to the response, preventing the leakage of private patient information. For example, researchers could query, "How many people in zip code ABCDE have diabetes?" and the differential privacy system would respond "12,045 people in zip code ABCDE have diabetes", which is a "blurry" response around the true value. If the query is too specific – e.g. "How many people in zip code ABCDE have Fields condition [an extremely rare disease]?", it might return that there are only one or two individuals with the disease, potentially leaking private information. To protect their privacy, a differentially private system would add noise and instead would return something like "Five people in ABCDE have Fields condition", which is quite different from the underlying reality.

Use in financial services:

This technique is sufficiently mature to be operationalized within financial institutions; the potential benefits are clear, and the incremental costs of integrating such techniques into existing data systems are not excessive. Adding noise directly creates a trade-off between precision and privacy, and thus the technique is best-suited to evaluating general trends, rather than anomaly detection (e.g. fraud analysis) or accurate pattern-matching (e.g. optical character recognition). Several companies such as Immuta have operationalized differential privacy solutions and serve financial institutions today.



Technique #2: Federated analysis

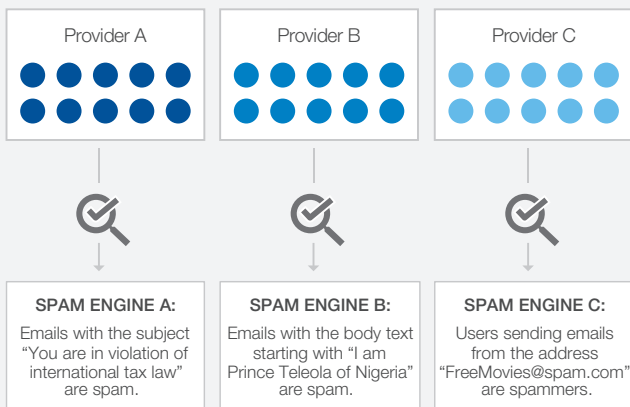
Overview:

If an institution is seeking to analyse large sets of data held across multiple databases or devices, it can combine them into one database to conduct analysis across the aggregate set of information. However, this introduces three issues. In some cases, the institution may not have permission to transfer the locally stored data (e.g. due to privacy or other localization restrictions in different jurisdictions). Furthermore, the data may be sensitive in nature (e.g. medical records, private transactions) and the data subjects (i.e. customers) may not feel comfortable sharing access to it. Finally, the centralization of data introduces a risk that if the central database is breached by a malicious third party, a gold mine of sensitive information would be exposed. As a result, both institutions and the data subjects themselves may be hesitant to share data in this way. One way to address these issues is to conduct the analysis on the disparate datasets separately, and then share back the insights from this analysis across the datasets.²³

In recent years, federated analysis has emerged as a solution to these issues, and the technique has been widely used by large technology companies (e.g. Google) to learn from user inputs on personal computing devices such as phones and laptops.²⁴ Research in this area is ongoing, and federated analysis models are being used in conjunction with other emerging technologies such as AI. For example, in March 2019, TensorFlow (a widely used open-source library for machine learning) published TensorFlow Federated,²⁵ an open-source framework that allows machine learning to be performed on federated datasets.

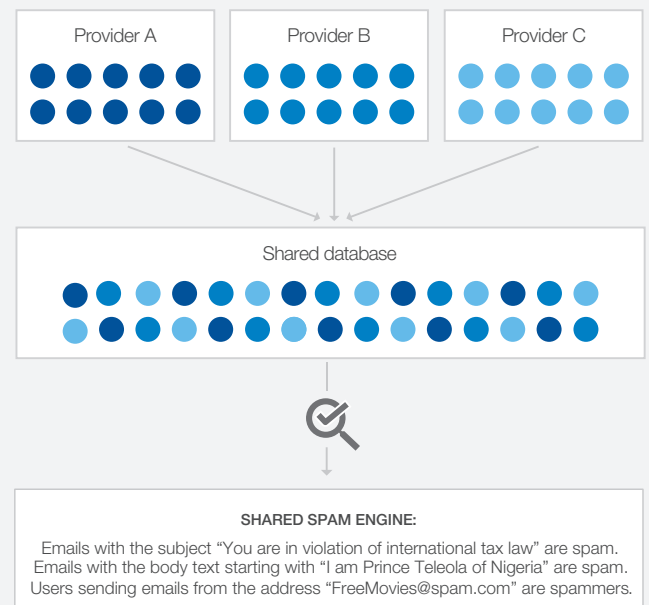
How it works:

Consider a hypothetical case where three email providers are seeking to reduce the amount of spam that their customers receive. One option would be to individually develop spam filters by analysing the emails that are reported as spam on their respective datasets.



In this case, the institutions would be duplicating their efforts as the characteristics of spammers are likely shared across each of their three customer bases. Furthermore, any differences in their analysis or input datasets would lead to gaps in their respective spam-detection engines.

To address these gaps, the institutions could instead combine their reported spam email data into a central database, and then create a shared spam-detection engine.

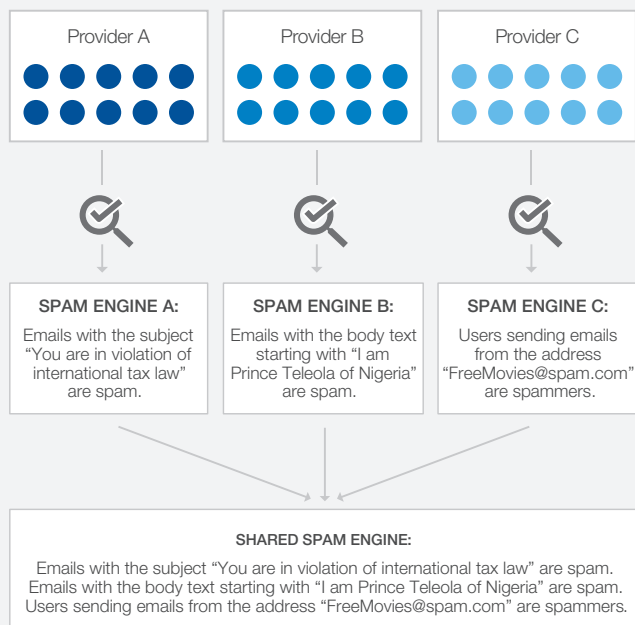


This engine benefits from the scale of data across the three institutions, leading to a superior product from which all customers can benefit. However, this solution introduces several issues: The customers of each email provider may not want their emails to be shared with third parties (even if the stated purpose is to improve spam filters for their own benefit). Furthermore, each institution has introduced the risk of exposing competitive information (e.g. who their customers are). Finally, this shared database presents a concentrated target for malicious third parties – breaching this single database would provide access to the sensitive information of customers across all three email providers. While this approach to data-sharing achieves the intended goal of an improved spam engine, it also introduces significant risks.

Instead, federated analysis can be used to achieve the same goal without introducing these new risks. Rather than sharing the underlying data, the institutions can share their spam-detection models and create an aggregated model.

This approach still results in a robust spam-prevention engine, while mitigating the risks that sharing the underlying data introduced. The institutions are able to benefit from a larger scale of data, while respecting any restrictions they may have on sharing customer data as that data is never shared with other email providers. From a security perspective, there is also no concentrated target for malicious third parties to attack.

It is important to note that this model does not necessarily produce an equivalent model to the one that would be derived by first combining the training data into a central location; in most cases, a model trained through federated machine learning would be inferior to the one trained on a centralized dataset. An example demonstrating this is shown in [Use case #1](#).



Where it could have helped:

In 2017, security researchers were able to access the personal data of 31 million users of an Android app called ai.type²⁶ – a third-party keyboard that allowed users to customize their phone/tablet keyboard and offered personalized typing suggestions. The app collected various types of data (e.g. contacts to offer those names as suggestions, or keystroke history to improve the auto-complete functionality) and stored this information in a single, central database. This database was then cleansed of private information (e.g. anything typed in password fields) before being analysed to provide autocomplete suggestions. However, researchers were able to access the database before this cleansing was performed, and were able to expose the email addresses, passwords and other sensitive information of all 31 million users. Rather than centralizing the data in one location, ai.type could have used federated analysis to create local predictive models on every user's phone. These could then have been aggregated across the app's 31 million users rather than the data itself, protecting the typing history of individual customers.²⁷ The aggregate model could then have been pushed back to individual phones in an update, and the learning process could have been continuously repeated; this would have allowed the keyboard to provide advanced recommendations based on its aggregate userbase. This is the approach that Google and Apple have taken with the default keyboards offered by Android and iOS.²⁸

Use in financial services:

While this technique is well-understood and mature from a technical perspective, its application in the financial services industry to date has been limited. The value of federated analysis is greatest when the number of separate sources of data is high – e.g. on cell phones, IoT (internet of things) devices, laptops, etc. Within financial services, rarely is sensitive information stored across this scale of hundreds of thousands of separate sources of data. Rather, transactions, customer information, etc. are stored centrally by the financial institution, and in most major geographies the top 10 players serve the majority of the market. However, federated analysis is a technically mature methodology and can still drive benefits in the financial services industry; one such use case is explored in [Chapter 3](#).



Technique #3: Homomorphic encryption

Overview:

In some cases, data analysis needs to be conducted by a third party, for one of two reasons:

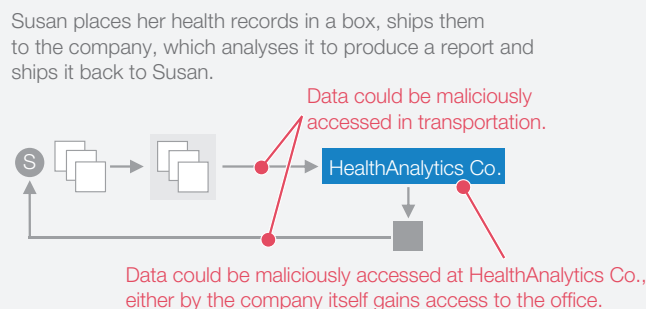
- The third party has capabilities the data steward does not, and the third party wishes to provide their analytics as a service without sharing the underlying functions they are using
- The third party has access to other, complementary data that the data steward does not have, and as a result is able to provide better analytics and insights than the steward could do independently.

As with federated analysis, however, the data steward may not have permission to transfer the data. Furthermore, if the data steward does not trust the third party, it will be reluctant to share this data for fear that it will be misused by insiders within the third party or its other partners. Finally, if this third party were to be breached, the original data steward would likely still be held responsible by its customers for sharing the data with the third party in the first place. Homomorphic encryption (HE) can be used to address these challenges by encrypting the data so that analysis can be performed on it, without the information itself ever being readable. The results of the analysis would also not be readable by anyone other than the intended party (usually the owner of the input data).

Homomorphic encryption was first theorized in 1978, accompanying the development of the “RSA” cryptosystem in 1977 – one of the first encryption schemes widely used to transmit data.²⁹ Under RSA, a (public) key is used to encrypt data and make it unreadable. This data can then be transported to the intended recipient, who decrypts it using a different (private) key. In 1978, the question was raised whether data could be encrypted in a way that would allow for different types of functions (e.g. addition, multiplication) to be performed without first decrypting the data and thus exposing sensitive information. For over 30 years, solutions were proposed that allowed for a specific function to be performed, but a fully homomorphic system where any transformation could be performed was not found. In 2009, the first fully homomorphic encryption (FHE) system was proposed by Craig Gentry³⁰ and throughout the 2010s, significant advancements were made in the efficiency and viability of FHE systems.

How it works:

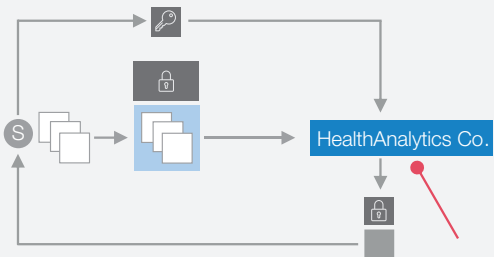
Consider a situation where Susan is looking to conduct sophisticated analysis on her health records to identify and predict any potential risks to her well-being. However, she doesn't herself have the capabilities to conduct such analysis and instead relies on a third party: HealthAnalytics Co., the leader in this field. To share her data with HealthAnalytics Co., Susan could collect all of her health records into a box and ship it to the company, but this introduces several risks: The box could be intercepted by an unauthorized third party (either in transit or once at HealthAnalytics Co.'s office); furthermore, malicious actors employed by HealthAnalytics Co. itself could use these documents for an unintended purpose.



Instead, Susan could use encryption to protect her information. In this case, Susan would collect all of her health records into a safe and ship it to HealthAnalytics Co. without the key and send the key to them separately through a different channel. This eliminates the risk of the contents of the safe being accessed by an unintended party: Even if the safe were to be accessed during transportation or at HealthAnalytics Co.'s office, a malicious third party wouldn't be able to open it without the key. A bad actor would have to breach both the HealthAnalytics Co. database and the transportation

channel that Susan uses to share her key to access the data, reducing the security risk. However, once given the key, Susan cannot be sure that the company itself will not use the documents for unintended purposes or make any copies of it. Thus, this form of “encryption” is also not completely secure.

Susan places her health records in a locked safe, ships it to the company and separately provides the key so that the data within the safe can be analysed. This analysis is placed in another locked safe and shipped back to Susan.



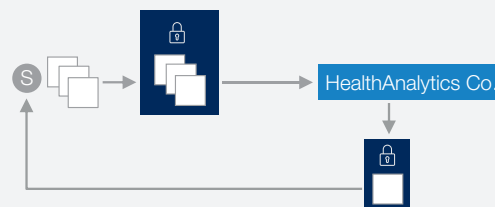
Data could be maliciously accessed at HealthAnalytics Co. by a bad actor within the company who has access to both the safe and the key, or an external bad actor who accesses the data during its analysis (when it has been removed from the safe).

To completely protect her data, Susan could use homomorphic encryption – which is effectively a special type of safe. She locks her health records into this special safe and sends it to HealthAnalytics Co., without the key. If a third party attempted to access the safe during

transportation or while it is in HealthAnalytics Co.’s office, they would not be able to (as they do not have the key). Unlike the previous case, this special safe allows HealthAnalytics Co. to conduct the required analysis on the safe itself, without ever opening it. The analysis on this special safe transforms it into another special safe containing the results, which can also be unlocked only by the key still held by Susan. HealthAnalytics Co. then ships this safe back to Susan, who uses her key to unlock it and read the analysis of her health records.

The company itself is not able to read the health records or even the results of the analysis it conducted, since they are protected by the special safe. Throughout the transportation/ storage of the information, it is also protected by the same key held by Susan.

Susan places her health records in a homomorphic encryption safe and ships it to the company. The company analyses the safe as if it were the underlying health records, producing another safe that can be unlocked only by Susan. The safe itself is shipped back to Susan, who uses their key to turn it into the underlying report.



Where it could have helped:

In 2018, the story broke about Cambridge Analytica, which had amassed data on more than 50 million Facebook users.³¹ The company purchased the data from a personality quiz app that collected users’ names, emails, profile photos, friend networks, likes and other information, and provided users with a high-level personality profile in return. The app stored the data it “scraped” and later shared it with a third party, Cambridge Analytica, which built detailed psychographic profiles to target audiences with digital advertisements. One possible approach to prevent such misuse of data (though probably not the most efficient or direct way of achieving this goal) would be to use homomorphic encryption – either mandated by Facebook or voluntarily used by the personality quiz app as a responsible data steward. With homomorphic encryption, users’ data would be encrypted before it was shared with the third-party personality quiz app. The app would then analyse this encrypted data and return a personality profile to individual users that the app itself cannot read. Users would be able to decrypt these results with their private key (based on their Facebook password), and the data itself would not be usable or even readable by Cambridge Analytica or any other third parties.

It is critical to note that encrypting the data (homomorphically or otherwise) does not free institutions from their privacy obligation. The data in this case would still fundamentally be personal information, and require robust data management and oversight to ensure it is shared and used in an ethical manner.

Use in financial services:

Generally, at the current level of sophistication, the use of homomorphic encryption at scale is limited for two key reasons: the limitations of the techniques and the lack of widely accepted standards.

Many homomorphic encryption schemes allow for only one type of operation (e.g. addition or multiplication, but not both), and analysis on data that is fully homomorphically encrypted (where any type of operation is possible) is several orders of magnitude slower than the same analysis on unencrypted data. As a result, the use of this technique is limited to use cases with a narrow set of functions (in the case of HE), or where the speed of calculation and cost of computation are not a priority (in the case of FHE). However, recent improvements in these techniques allow for some computations to be completed in relatively short order (seconds and minutes), enabling applications of homomorphic encryption to protect highly sensitive data. This remains an area of active development, and start-ups such as Ziroh Labs and Inpher have developed HE and FHE schemes that are computationally viable for real-world use cases.

Legacy encryption systems have widely accepted standards that allow for a high degree of interoperability and widespread use. No such widely accepted standard exists for HE or FHE schemes, greatly diminishing the usability of any given homomorphic encryption scheme. There are some initiatives underway (e.g. Homomorphic Encryption Standardization) that seek to define community standards for this technology.



Technique #4: Zero-knowledge proofs

Overview:

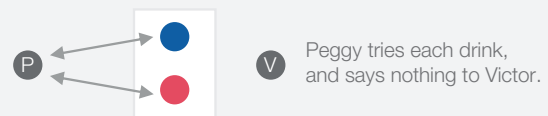
Sometimes, users seek to share specific information without leaking any additional data. This is important in situations where the user seeking to share the information does not trust the other party not to use it for something other than the intended purpose. For example, when filling out a rental application, an individual may want to prove that their income exceeds the landlord's minimum requirements. However, they may not want to share exactly how much they earn – if it is significantly over the minimum requirement, there is a risk that the landlord will raise the rent at the first available opportunity. In this case, the third party receiving the income verification could use the additional information they received (the exact salary) to derive additional knowledge that the applicant sought to keep secret. Zero-knowledge proofs (ZKPs) allow for one party to prove to another some specific information without sharing anything other than the intended information.

ZKPs were first introduced in 1985 in the paper “The Knowledge Complexity of Interactive Proof-Systems” by Shafi Goldwasser (MIT), Silvio Micali (MIT) and Charles Rackoff (University of Toronto).³² Since then, ZKP logic has continued to evolve to include a broader set of use cases, including witness-indistinguishable proofs, non-interactive proofs, quantum-resilient proofs and more. As with federated analysis, the technique is being used in conjunction with other emerging technologies – most notably with distributed ledgers to enable the transfer of assets across a P2P system with complete privacy.

How it works:

Consider the hypothetical situation where Peggy wants to prove to Victor that she can tell the difference between two types of soda, stored in two identical glasses. Peggy has two additional desires: She wants to keep her method for distinguishing between the two (say, by knowing that one is sweeter) a secret from Victor, and she does not want to let Victor know which glass is which brand of soda. If she is able to do this, she would have “zero-knowledge-proved” that she can tell the difference between the two drinks, without exposing any other information about herself or the contents of the glasses.

In order to do this, Peggy should sample each glass, then turn away from the table. Victor should then randomly either switch the glasses or leave them in the same position (with approximately 50% probability of doing either), then allow Peggy to sample each glass again. Peggy should respond by stating whether the glasses were switched or not, but should not communicate which glass contains which brand of soda, or how she knows if the glass was switched or not. The first time this test is conducted, Peggy has a 50% chance of being right just by guessing. However, if she truly can tell the difference between the two, she will be able to consistently answer correctly as the test is repeated, and the chances of her guessing the right answer decreases significantly.



By the 20th trial, there is an approximately 1/1,000,000 chance that Peggy is guessing, and thus Victor can be reasonably certain that she knows the difference between the two soda brands. This proof is zero-knowledge, as Victor does not know which glass is which, and also does not know how to differentiate between the two sodas.

Where it could have helped:

In January 2019, an employee at a major American retailer was arrested for allegedly sharing customers' credit card numbers with an accomplice who would then make fraudulent purchases using the stolen card information.³³ The employee would memorize and transcribe customers' card numbers while ringing through their purchases, and text the numbers to the accomplice shortly after. The accomplice would then use the stolen card information to purchase gift cards, sometimes giving the employee gift cards for her alleged role in the theft. Similar credit card theft schemes are responsible for a share of the estimated \$130 billion in card-not-present fraud that retailers are expected to encounter between 2018 and 2023.³⁴ We can now envisage how a zero-knowledge proof payment system could prevent such losses by allowing individuals to validate their bank information and balances at a retailer without ever exposing their account information and CVV code to any third party (e.g. the cashier).

Use in financial services:

ZKP has only recently seen real-world operational uses as the methodology continues to mature, but it has applications across a variety of use cases – including payments (e.g. Zcash³⁵), internet infrastructure (e.g. NuCypher³⁶), digital identity (e.g. Nuggets³⁷) and others. Large institutions such as ING have invested in advancing ZKP techniques in financial services,³⁸ and it is expected to be a critical enabler of distributed ledger technologies more broadly (as it allows individuals and institutions to protect private information on public distributed ledgers).



Technique #5: Secure multiparty computation

Overview:

As with homomorphic encryption and zero-knowledge proofs, this technique allows for individual privacy to be maintained when sharing information with untrusted third parties. Secure multiparty computation (SMC) allows institutions to conduct analysis on private data held by multiple other institutions without ever revealing those inputs. In the past, doing this would have required an intermediary to act as a middle man to the data-sharing, which however introduces several issues:

- Insiders within this intermediary could misuse the data (e.g. sell it to another party seeking to use it for an unintended purpose). Within the context of collaborative endeavours, the third parties/intermediary may even be competitors (e.g. banks sharing transaction data to identify payments fraud), which raises the risk that competitive secrets would be exposed.
- If the intermediary were breached by an external bad actor, institutions' sensitive data would be exposed, and the institution would likely still be held responsible by its customers and regulators, despite not being directly responsible for the security breach.

With SMC, the intermediary is replaced by an incorruptible algorithm that, even if breached, does not expose any sensitive information. Fundamentally, SMC relies on "secret sharing",³⁹ where sensitive data from each contributor is distributed across every other contributor as encrypted "shares". These shares, if intercepted by a malicious third party or misused by an individual contributor, would be worthless, since they are decipherable only once they are combined with the information distributed across many other parties.

In the late 1970s, as computing became common in homes and offices around the world, SMC first emerged as a solution to the problem of establishing trustworthy systems in environments with no trusted third party (e.g. how can I play poker online when I cannot trust that the website running the game is not rigging the system?)⁴⁰. Schemes developed since then have evolved over time to address a broader set of use cases, and the first live implementation of SMC was in 2008, when it was used to determine sugar beet market prices in Denmark without revealing individual farmers' economic position.⁴¹ Since the early 2010s, research has focused on improving the operational efficiency/scalability of SMC protocols.

How it works:

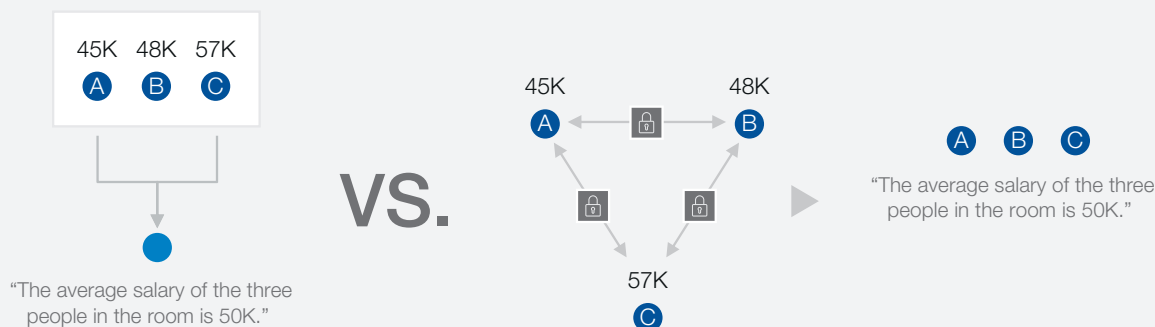
The specific logic underpinning secure multiparty computation is particularly complex, perhaps more so than the other techniques outlined in this paper. In the interest of ensuring that the fundamental process is understandable without significant technical expertise, we have provided two varying descriptions of the technique:

1. High-level summary: A short, abstract and high-level outline of the technique and its benefits.
2. Detailed explanation: A detailed case study that steps through a hypothetical explored in Technique #1 (differential privacy), with example calculations along each step.

Alternatively, a useful video by Boston University explaining SMC can be found [here](#).

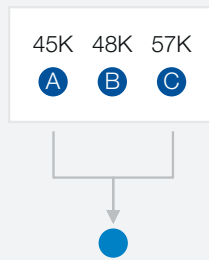
1. High-level summary:

Fundamentally, SMC relies on the sharing of encrypted messages among several parties, configured in such a way that through the required analysis and calculations, sensitive data is not shared between parties but the correct end result can still be derived. SMC systems can be configured in such a way that each party is responsible for a portion of the calculation, so there is no need for a trusted intermediary.



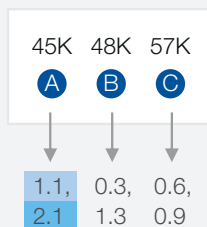
2. Detailed overview:

Let us return to the hypothetical example discussed in Technique #1, simplified with three individuals in the room instead of 10. In the original example with a trusted intermediary, the process of knowing the average salary of the room is relatively simple.

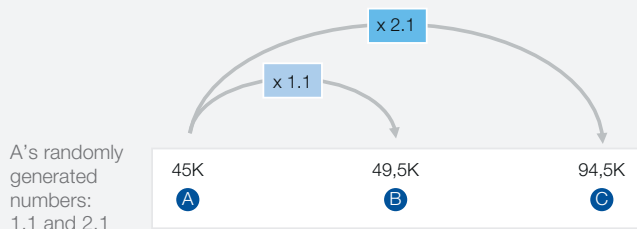


"The average salary of the three people in the room is 50K."

This hypothetical relies on an assumption that that this intermediary is trustworthy and incorruptible, something that cannot always be taken for granted. It is possible that the intermediary is colluding with one of the individuals in the room (or a third party) and later shares the private information; or that the intermediary's records are breached, and a third party can access the private information without consent. SMC can be used to mitigate these risks – rather than involving an intermediary, an algorithm can be used to perform the same function. To start, each party randomly selects two numbers between 0 and 3 (the upper limit being the number of participants in the data-sharing collaboration).



Each participant then multiplies their salary figure by the randomly generated figures and provides the other two participants with this distorted number. Let us walk through the maths that participant A performs:



At this point, B and C each have a copy of A's salary that is wildly different from each other and cannot be reverse-engineered into the original number even if the two were colluding. Participants B and C also perform the same exercise, using their own randomly generated numbers to modify the salary figures they share with the other participants. This creates a matrix of warped salary information.

	A	B	C
A provides ...		49.5K	94.5K
B provides ...	14.4K		62.4K
C provides ...	34.2K	51.3K	

In order to devise the total salary of the group (to then divide it by three and calculate the average salary of the three individuals), each participant takes their private salary information, adds the numbers that were provided by the other participants (the column with their name) and subtracts the numbers that they provided to the other participants (the row with their name). Let us walk through A's process:

	A	B	C
A provides ...		49.5K	94.5K
B provides ...	14.4K		62.4K
C provides ...	34.2K	51.3K	

Participant A would add the blue highlighted numbers to their actual salary figure and subtract the orange - actual salary figure and subtract the orange - highlighted numbers from their salary figure. A then shares this result with the other participants.

$$\text{Thus, A's response would be} = 45 + 14.4 + 34.2 - 49.5 - 94.5 = -50.4K$$

Participant B executes the same process with his/her own salary figure:

	A	B	C
A provides ...		49.5K	94.5K
B provides ...	14.4K		62.4K
C provides ...	34.2K	51.3K	

$$\text{B's response would be} = 48 + 49.5 + 51.3 - 14.4 - 62.4 = 72K$$

Participant C does the same:

	A	B	C
A provides ...		49.5K	94.5K
B provides ...	14.4K		62.4K
C provides ...	34.2K	51.3K	

$$\text{C's response would be} = 57 + 94.5 + 62.4 - 34.2 - 51.3 = 128.4K$$

Added together, these three responses equal the total salary of the three individuals in the room, which can then simply be divided by three to derive the average salary:

$$50.4 + 72 + 128.4 = 150K$$
$$150K/3 = 50K$$

It is important to note that at no point during the entire process were any of the participants' actual salary figures revealed: 45K, 48K and 57K were not seen in any of the intermediate steps. Nor is it possible to reverse-engineer those figures from any of the intermediate inputs provided by the participants, since those inputs were warped by the random modifiers (numbers between 0 and 3).

However, since all parties learn the true and exact output from the analysis, one party may still be able to cross-reference the output with other information in order to infer some sensitive data (as seen in the case for Technique #1, where one party is able to reverse-engineer an individual's salary information by deducing it from the average salary and known salary figures of other participants who contributed to that average). Differential privacy can be applied to the outputs of an SMC system to provide privacy not just in the analysis of the data, but in the sharing of the results of the analysis as well. This is explored in greater detail through the use cases in the following section, where we explore how different techniques can be combined in real-world applications in financial services.

Where it could have helped:

On 10 February 2009, the US's Iridium 33 communications satellite collided with a Russian Kosmos 2251 satellite, instantly destroying both.⁴² The positional data on board each satellite, if shared between the United States and Russia, could have detected and prevented the impending collision, but satellites' orbital data are guarded very carefully for the national security and privacy of both the citizens and the military of each country. An SMC protocol could be used to enable the sharing of only the key insights (i.e. "Will any of the United States' and Russia's respective satellites collide in the near future?") without sharing the underlying location data.

Use in financial services:

SMC is a relatively nascent technique, and as a result its application in the financial services industry (and more broadly) is limited. This is in part because SMC requires a completely customized set-up for each use case, creating extremely high set-up costs (unlike, for example, differential privacy, where generic algorithms can be used across use cases). However, "compilers" that abstract the underlying protocols to enable general-purpose computing are being developed, supporting data science and machine-learning applications more broadly.

Current SMC systems have high communications costs, making them expensive to operate on an ongoing basis. The technique is continuing to evolve, though, and fintechs such as Inpher (with an investment from JPMorgan) have developed SMC products and services specific to the financial services industry. This technique's use in the industry will likely continue to grow over time.

Chapter 3: Applications in financial services

Each of these techniques has different advantages and disadvantages, with a host of potential applications across the financial sector. It is important to note that these techniques do not need to be applied exclusively; in fact, combining them may enable highly targeted mixes of privacy, security and utility, with the benefits of one technique being used to reinforce the limitations of another. In this chapter, we explore how privacy enhancing techniques can be used to navigate privacy tensions and enable institutions to unlock new value by sharing data in new ways.



Unlocking new value for financial institutions

Use case 1: Detecting vehicle insurance fraud

With federating learning, differential privacy and zero-knowledge proofs

Context:

In the United States, the total cost of non-health insurance fraud is estimated to be more than \$40 billion per year, costing the average family between \$400 and \$700 in increased premiums annually.⁴³ In the vehicle insurance industry specifically, the cost of fraud is shared between customers (who pay higher insurance premiums than would otherwise be needed to insure the actual risk) and financial institutions (which make payments on fraudulent claims that eat into their loss ratio, and thus ultimately their profitability).

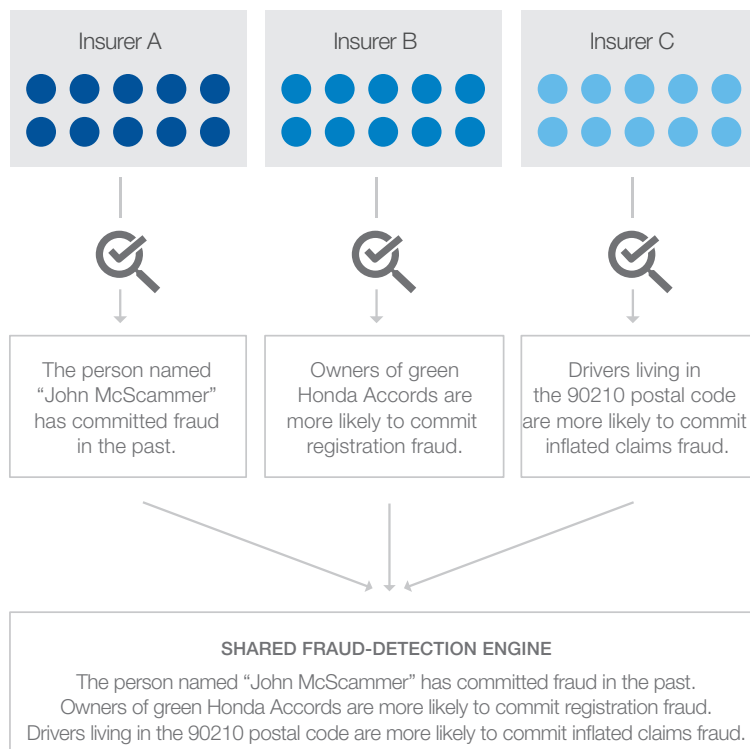
Data-sharing opportunity:

There is an opportunity for insurers to share data in order to reduce fraud, using registration, claims, telematics, insured assets and customer data, as well as other unstructured data such as medical reports. This would allow institutions to realize two benefits:

- An increased scale of shared data, leading to improved predictions and analysis. For example, an increased scale of claims data, telematics data and other unstructured data would allow institutions to better identify patterns that suggest fraudulent claims.
- The identification of duplicate claims, filed against the same assets or the same incident across multiple insurers.

However, much of this data is sensitive and accompanied by significant privacy concerns. Customers would not want their private information such as their registration data, claims data, personal information and other data such as medical reports being shared with third parties. Insurers themselves would be wary of sharing such information with their competitors, since it could be misused to deduce underwriting and pricing strategies, and other sensitive, competitive information.

Hypothetical application of privacy enhancing techniques: Federated analysis could be used to create master fraud detection/prevention models across registration and claims, without ever sharing the underlying customer data across institutional lines. This would allow insurers to benefit from an increased scale of data, while protecting the privacy of their customers and the confidentiality of their business operations.



That said, it is important to note that federated analysis may not capture the complete value of duplicate information across institutions. For example, consider a situation where an individual named John McScammer has vehicle liability insurance coverage with Insurer A and health insurance with Insurer B. Mr McScammer then gets into an accident and files a claim with both Insurer A and Insurer B, getting paid twice. Federated analysis would miss this “double dipping” since the individual claims filed and analysed within each insurer’s isolated dataset would be legitimate.

A variety of other privacy enhancing techniques could be used to address this gap in analysis, depending on the exact architecture of the data-sharing collaboration. Insurers could combine their datasets into one central (homomorphically encrypted) database on which to conduct analysis, which would be able to catch duplicate claims as in the case of Mr McScammer above. This central database could also be queried/analysed directly to arrive at the same insights that a federated analysis model could deliver, using differential privacy to ensure that the privacy of individual customers does not leak through the analysis. ZKPs could also be used to query each insurers’ individual datasets as part of the claims process. For example, when Mr McScammer files a claim at Insurer A, they would query B and C to see if claims have been filed by Mr McScammer, or on the same underlying insured assets in recent history. By using ZKPs, these queries would be able to check for a match without exposing the specific customer or assets in question to Insurer B or C, preventing a leakage of the sensitive information (e.g. the name of the customer from Insurer A to B or C).

Use case 2: Becoming the trusted guardian of data With zero-knowledge proofs

Context:

Historically, technology firms with modern systems, sophisticated analytics and proprietary datasets have been able to establish themselves as stewards of customers’ sensitive data (e.g. email) and identity (e.g. social sign-on). However, in the wake of numerous data privacy scandals and fines issued under regulations such as GDPR, the limitations of this practice are being brought to light – namely, that many technology firms’ core business model relies on monetizing the data they are provided.

Data-sharing opportunity:

Financial institutions, unlike technology firms and many firms in other sectors, have not historically relied on monetizing data as a source of revenue. Financial institutions are also already subject to stringent regulations in regard to the security of the data they hold, and have developed brands over many decades and centuries as trusted custodians of another valuable asset – money. As such, this positions financial institutions as the drivers of the next generation of data stewardship in society: to assert a new model for digital services based on trust and regulatory obligation. This presents financial institutions with the opportunity to deepen their engagement with customers through more frequent interactions. It also offers the chance for institutions to expand their presence in financial services-adjacent products and services.

Hypothetical application of privacy enhancing techniques:

To illustrate how financial institutions can serve as the trusted guardian of data, let us return to the hypothetical case briefly introduced in the overview of zero-knowledge proofs: an individual seeking to prove to their landlord that they exceed the landlord’s minimum requirements, without exposing their specific income (or in some cases, specific employer). As the recipient of the individual’s direct-deposit pay cheques, a retail bank would know this information already and is positioned as a trusted authority (i.e. one that a reasonable landlord can trust).

Privacy enhancing techniques are not necessarily needed to satisfy the landlord: A notarized letter on bank letterhead would likely suffice and is indeed used in the current state in some rental markets. However, the use of zero-knowledge proofs would provide two specific benefits:

- It would allow the customer to self-serve the necessary documentation, rather than relying on a financial adviser or branch customer service representative to assist in creating a notarized letter; this allows for faster service at lower costs
- It would be more trustworthy, as a ZKP system could be easily verified by the landlord directly, whereas a notarized letter could be forged or edited by the applicant and would be difficult for the landlord to validate.

It would not be worth the investment for any financial institution to undertake the significant technology spend required to build a system just for ZKP income verification. However, a similar mechanism could easily be expanded to many more customer attributes. This could include financial attributes such as transaction data, as well as non-financial attributes such as age and address. Indeed, it is highly likely that banks have more up-to-date information on many individuals than the traditional sources of identity verification (such as passports or drivers’ licences), since banks are required by law to keep such information up to date, while the same updates are relayed to government services only at passport/licence renewal every few years.

By collaborating, financial institutions would be able to unlock even greater value as trusted custodians of customer data: Any one institution will hold several pieces of data (e.g. debit and credit card transactions), but likely not all pieces of data (e.g. mortgage balance, investment balances). A collaborative network of data stewards would be able to route incoming requests from third parties to the appropriate financial institution on a case-by-case and customer-by-customer basis.

Hypothetically, this could allow customers to do the following and more:

- Validate their age without disclosing it specifically (e.g. to rent a car without paying a youth tax)
- Authenticate into government services (e.g. tax assessment statements) and financial services (e.g. free credit score checks) easily and instantly
- Share their credit score within the specific bands of a lender’s decision-making system’s ranges, without sharing the exact score

Use case 3: Mimicking open banking without regulation

With secure multiparty computation

Context:

Open banking regulation has been observed in multiple jurisdictions around the world, including the UK, EU, Japan, Australia, Hong Kong and Canada. However, in other parts of the world, it has not yet become a regulatory requirement, and certain environmental characteristics make it difficult for top-down regulation to be effective. For example, the United States banking landscape is heavily fragmented, both from a regulatory perspective (where banks are regulated both on the state level and federal level) and from an institutional perspective (with over 5,000 FDIC-insured institutions in the country).⁴⁴

This fragmentation makes it difficult for any unifying regulation, as observed in other jurisdictions, to mandate what data institutions are required to share and how. Many institutions perceive this to be an advantage: They do not need to make data available to third parties that would in most ways serve as competition to themselves, or undergo the extensive technology spend to make this data shareable.

Data-sharing opportunity:

However, this data is still being accessed by third parties through “screen-scraping” services. These services ask for a client’s online banking username and password, then use an automated system to log in and extract the client’s transactions on a regular basis. This introduces significant security risks (as clients are asked to share their credentials) as well as significant bandwidth usage (as the banks’ online web page is continuously loaded for the automated system to extract data, rather than just the relevant numerical and text-based information that is actually needed).

Furthermore, some institutions have realized that sharing data can also serve as a competitive weapon: It allows institutions to offer their customers value-added products and services. Building a data-sharing ecosystem defined by the institution rather than regulation allows for greater flexibility on what is in scope and the terms of the data-sharing agreements with third parties, allowing for greater control over the ecosystem. For example, BBVA has built a “Banking as a Service” platform that allows third parties to verify customer identities, move money and even originate accounts through code.⁴⁵ This has shifted BBVA from a competitor to a service provider for fintechs, and ultimately a holder of their deposits.

Hypothetical application of privacy enhancing techniques:

Secure multiparty computation can be used by financial institutions (assuming the technical challenges discussed previously can be overcome) to ensure that the data they share with third parties (e.g. customers’ transaction data) is used only for the intended purpose. For example, consider an institution partnering with a fintech to provide cash-flow forecasting for its small business clients; the fintech connects in real time with customers’ invoicing software (e.g. Xero) as well as the bank to automatically build detailed cash-flow forecasts and predict where bridging loans may be necessary.

While there is value being generated for the customer through this relationship (and potentially for the bank’s

lending business), the institution has a responsibility to ensure that their customers’ data is not misused. Under a typical data-sharing agreement, or one powered by screen scraping, the fintech would have full access to the customer’s transaction data and could potentially misuse this data without the bank or the customer knowing.

By defining a secure multiparty computation system, only the analysis that the system was designed for would be possible (e.g. summing up transactions, identifying recurring inflows and outflows), and the fintech would not be able to access the transaction data to the line item. This reduces the risk of abuse by the third party, or the potential for sensitive information to be leaked if the third party were to be breached. This allows for greater trust from both the institution entering the partnership and potential customers using the fintech’s service.⁴⁶



Unlocking new value for customers

Use case 4: Using an intelligent automated PFM adviser

With differential privacy

Context:

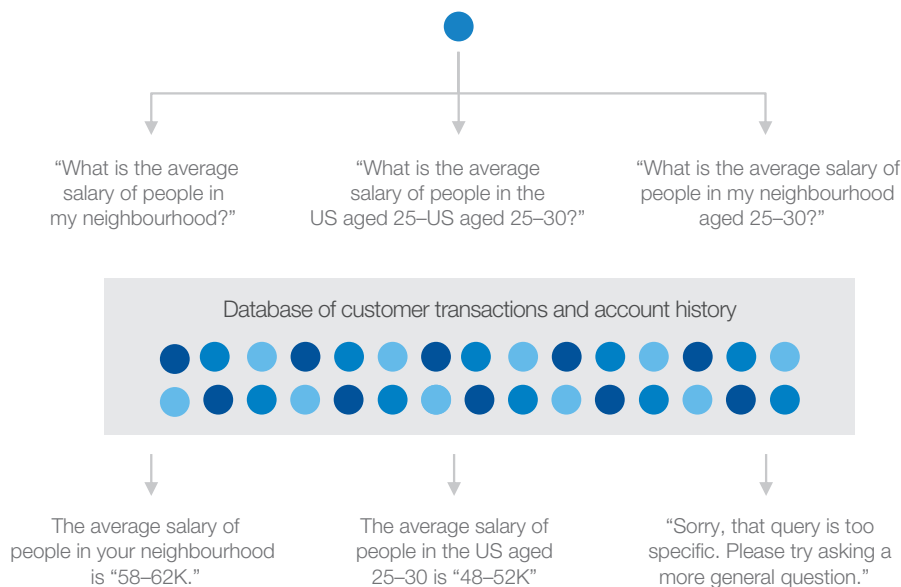
Privacy enhancing techniques can be used to enable competitive processes as well as collaborative ones. With increased access to data (through open banking) and the increased sophistication of automated analysis, a range of actors are becoming more focused on developing personal finance managers (PFMs) for the mass market. The potential benefits of such a product are clear – only 30% of Americans have a long-term financial plan that includes savings and investments,⁴⁷ and almost half are “concerned, anxious or fearful about their current financial well-being”.⁴⁸

Data-sharing opportunity:

Automated analysis across an entire database of retail banking activity could be used to provide sophisticated “people like you” recommendations. Open banking makes this data available to third parties (with customer consent), who can then provide personalized advice based on an understanding of the aggregate customer base; for example, an institution could answer customer queries such as “How much more/less than the average person of my demographic do I spend at bars?” While insightful, some customers may not feel comfortable with their spending habits being shared, even anonymously, with other users. If the demographic pool is small enough, the “people like you” comparison may be small enough for individuals to learn about the spending habits of specific other individuals that are also customers of the same PFM tool.

Hypothetical application of privacy enhancing techniques:

For the recipients of open banking data (e.g. an automated PFM adviser as described above), differential privacy can be a critical tool in unlocking the value of the cross-institutional dataset they aggregate. Differential privacy can be used to introduce noise into the process of generating insights and ensure that the privacy of the individuals in the dataset is not breached. This would break the privacy trade-off, allowing individuals to benefit from personalized and specific financial advice, while protecting the individual privacy of customers whose data informed it.



For larger financial institutions, there could be significant value in enabling such comparisons and analysis across jurisdictions, drawing on the scale of data available in one market to provide high-quality products and services in another. Imagine a large US bank with a successful automated PFM adviser that is leveraging differential privacy. To expand in Canada, this bank seeks to offer the same quality of advice from Day 1 to Canadian consumers who share many of the same consumer behaviours and preferences.

Due to privacy regulation, the bank may not be able to share transaction data, account balances and demographic customer information across borders, and as a result must acquire a large set of customers in Canada on which it can then perform the same analysis it has already conducted in the US. However, the use of differential privacy allows the Canadian institution to conduct and draw on analysis from US customers, without ever accessing the underlying data. This would allow Canadian customers signing up to benefit from the history of learning the institution has already established in the US, ultimately allowing the firm to address its cold-start problem: providing high-quality insights to customers from whom it has not yet gathered enough information directly.

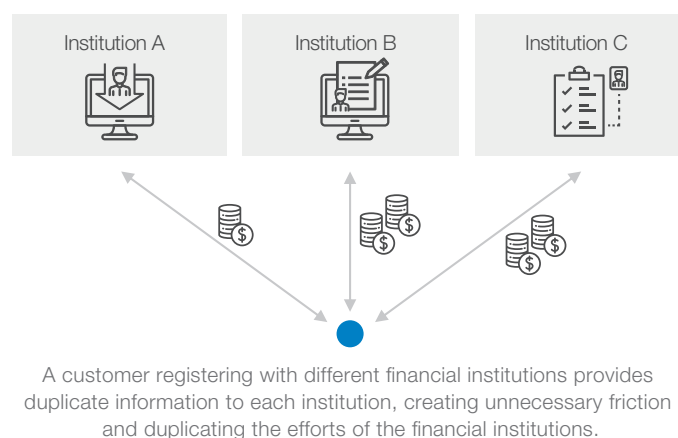
Use case 5: standardizing the customer registration programme for retail banks

With zero-knowledge proofs

Context:

Currently, retail banks manage their know-your-customer (KYC) and anti-money laundering (AML) onboarding processes independently. While there are many steps to the process, we will focus on the Customer Identification Program (CIP) for the purposes of this illustration. CIP is a requirement in the US for financial institutions to verify the identity of an individual wishing to conduct financial transactions through their infrastructure. At a minimum, this includes acquiring and verifying a name, date of birth, address and valid identification number at onboarding. For customers, this means submitting the same documents repeatedly when applying for products across financial institutions. For the financial services industry at large, this

creates a significant duplication of effort across institutions. While CIP is a US-specific obligation, institutions around the globe face similar requirements.

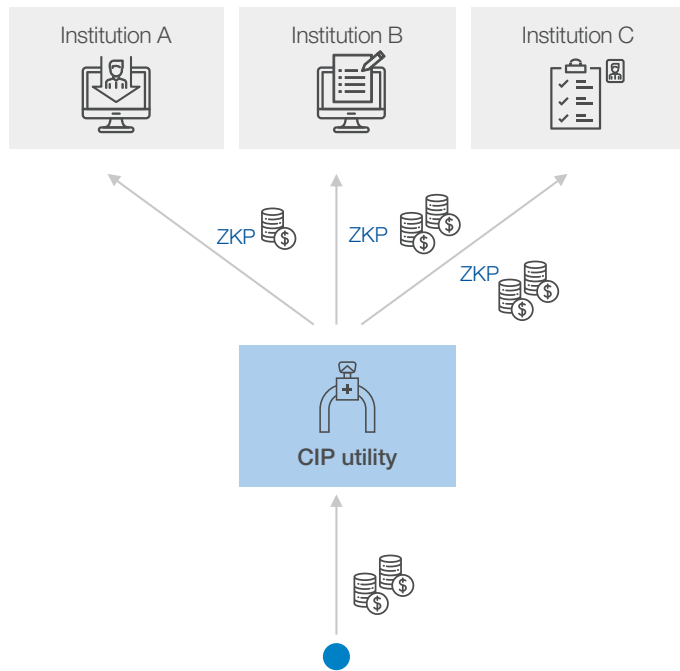


While burdensome, such processes are not perfect: In mid-2014, known fraudster Daniel Fernandes Rojo Filho was able to open 17 bank accounts at large financial institutions under his own name,⁴⁹ even though Googling it would have immediately revealed his history of financial crime.

Data-sharing opportunity:

By mutualizing the CIP onboarding process, customers would benefit from faster, standardized onboarding experiences with limited need for duplicate data entry across institutions. This would also enable more efficient KYC and AML processes more broadly, as such checks would be conducted against a single set of data. However, customer privacy regulation often prevents the sharing of personally identifiable information across institutions (and in some cases, across business units). Furthermore, institutions may prefer to keep their customer databases private out of fear of targeted competitive advertising campaigns on their customer base.

Hypothetical application of privacy enhancing techniques: Zero-knowledge proofs (ZKPs) can be used to address these concerns while reducing duplicate effort across institutions. Rather than each bank running its own CIP, a shared utility can be used to provide the necessary documentation. Users would sign up with this intermediary only once, providing all the traditional KYC documentation (e.g. address, identification number). Using ZKP, this intermediary would share the requested CIP data with different financial institutions as needed, providing only the information required for the specific product for which the individual is registering.



Through this system, the individual financial institutions would never need to store copies of an individual's identity data, preventing the creation of duplicate (and eventually outdated) information. Such a utility would be able to provide benefits beyond onboarding: In addition to simplifying the CIP experience, it could also support the ongoing monitoring processes required by KYC/AML regulations – e.g. a change in address would need to be provided only once to the CIP utility and would instantly be applicable to the individuals' various financial institution relationships the next time it is requested or needed by those firms.

Note that we are not proposing that the formation of a single, centralized utility is the ideal architecture for an end-to-end digital identity solution. Rather, we use this example to demonstrate how ZKPs could be used to address the creation of duplicate information as a result of CIP. For a more detailed read on the need for digital identity in financial services and other plausible architectures of digital identity solutions in financial services, please see the World Economic Forum's 2016 report on a blueprint for digital identity.⁵⁰



Unlocking new value for regulators

Use case 6: Analysing ecosystem-wide financial risk exposures

With secure multiparty computation or homomorphic encryption

Context:

Preventing systemic risks (which threaten entire financial systems and markets) from materializing is a complex task and is difficult to do ex ante. Institutions individually manage a large variety of risks (e.g. credit, liquidity) according to the regulatory requirements prescribed to them, but do not have all the information needed to see how those risks may be forming across the entire ecosystem; while their individual processes may be robust, their interactions with other market participants can lead to unexpected outcomes for the financial system as a whole. It is difficult for any individual institution or regulator to predict or proactively detect these risks, as the data required to "see the full picture" is fragmented across multiple bodies.

Data-sharing opportunity:

Proactive analysis of ecosystem-wide data has the potential to provide an advance warning of systemic risks being created across the financial system – such as those that led to the 2008 recession. For example, aggregating data across the US mutual fund industry could have highlighted open-ended funds' concentrated exposure to Lehman bonds. However, it is clear why such data is not collected in today's financial system: It can be highly sensitive and sharing it openly would pose significant competitive threats to institutions' strategies. At the same time, timely access to this information is critical to anticipating threats to the safety and soundness of the financial ecosystem.

Hypothetical application of privacy enhancing techniques:

As outlined in greater detail by Emmanuel A. Abbe, Amir E. Khandani and Andrew W. Lo in their paper "Privacy-Preserving Methods for Sharing Financial Risk Exposures",⁵¹ secure multiparty computation can be used to conduct the relevant aggregate analysis on financial institutions' risk exposures without breaching their individual confidentiality and revealing their strategies to competitors.

Using mechanisms like those outlined in [Technique #5](#) (secure multiparty computation), the aggregate values of loans by sector vertical (e.g. housing vs. industrial vs. vehicle) can be calculated to indicate the economy's sensitivity to changes in interest rates without exposing firms' individual credit portfolios, which is sensitive and proprietary data.

In theory, homomorphic encryption could also be used to conduct similar analysis, but the practical limitations of the technique on more complex analysis (e.g. means, variances) mean that it is probably too computationally expensive to provide meaningful insights in a timely manner. As homomorphic encryption techniques continue to mature, they may more directly substitute SMC systems.

Closing comments

To date, successful financial institutions have competed on the basis of price (i.e. offering products and services at the lowest total cost) and customer experience (i.e. offering unique value propositions). These bases of competition – aided by technological change – have driven the bulk of the advancement in the global financial system over the past several centuries. However, another pillar is emerging as a vital characteristic of winning institutions: privacy and security. For customers and regulators, knowing that an institution will safely store and manage data is critical to garnering trust, and in the wake of multiple scandals across several industries, this trust has been severely shaken.

As control over data increasingly shifts into the hands of customers (driven by customer demands and by regulatory mandates), there is a growing view in the financial sector that institutions will lose the ability to exploit the data they hold to create value: for customers, for themselves and for societies at large.

Understanding the opportunity of PETs

However, as demonstrated through the use cases explored in this paper, the emerging set of techniques known as “privacy enhancing techniques” (PETs) have the potential to create value that at first glance would be impossible to capture due to concerns about data privacy. A combination of differential privacy, federated analysis, homomorphic encryption, zero-knowledge proofs and secure multiparty computation can enable many uses outside of those mentioned already, including:

- Preventing insider trading by sharing patterns and insights from trade data across institutions without sharing the underlying trade data itself.
- Preventing bid rigging by replacing intermediaries with autonomous, transparent and incorruptible algorithms that perform the same service.
- Detecting tax fraud by analysing companies’ purchase and sales invoices while maintaining the data confidentiality of those transactions.

Challenges to realizing the opportunity of PETs

The opportunity presented by PETs is large and growing rapidly, but it is critical to *note that using privacy enhancing techniques successfully will require institutions to take several steps beyond understanding and deploying the techniques themselves.*

Investing in research and development: Many of these techniques are relatively nascent, with significant developments occurring over the past few years that have brought them into the realm of possibility for deployment within financial services. However, institutions will need to invest heavily in making these techniques easier to use in business applications. The bulk of development of privacy enhancing techniques to date has been based in academic research, with lesser consideration for implementation in financial services by developers and usage by business users. As a result, many of the systems today are difficult to translate to business contexts. Several companies have emerged to bridge this gap as a service, enabling

institutions to more easily take advantage of the benefits offered by PETs. Whether by collaborating with such third parties or by funding new research and development, institutions will need to invest in continuing the wave of innovation in PETs that has been observed over the past few years.

Collaborating with the public sector: Due to their nascence, there is uncertainty in some cases on how PETs would be treated under privacy regulations around the world. For example, federated analysis or secure multiparty computation in theory should allow institutions to analyse their data across regions where sharing data across international borders would otherwise be prohibited. However, ensuring that this is explicitly permitted by regulation would be important to preventing any fines or other regulatory risks from materializing, and in many cases the required regulatory certainty does not yet exist. Soliciting this certainty will necessitate an increased understanding of PETs as well as open discussions between the public and private sector on what is a safe approach to using these techniques in the financial sector.

Educating customers: Many of these techniques are unintuitive – and as a result risk creating experiences that do not feel private and secure, even if they are. In order to garner trust and adoption by customers, institutions will need to take a two-pronged approach to implementing PETs within their businesses – focusing both on protecting customers’ data and on helping customers *feel* protected.

Tackling related obstacles: Beyond the issues directly related to PETs, institutions will need to navigate several adjacent issues in order to fully realize the opportunities discussed in this paper and more. These challenges include:

- Poor data quality: historical datasets contain errors from manual entry, lack detail and/or are difficult to cleanse and format for computer processing.
- Legacy technology: ageing core systems do not support the type of data access (e.g. real-time, API-based) required to enable seamless data-sharing and analysis.
- Fragmented data architecture: Data within an organization is hosted across a variety of databases that cannot easily be integrated to generate insights.

- Lack of data interoperability: A lack of shared data formats often leads to the loss of depth or quality when sharing data across institutions.
- Geographic discrepancies: Different constraints and allowances on how data can be used across jurisdictions due to policy differences adds further complexity for global organizations.

For many of these adjacent challenges, various technologies are emerging to help address these issues. For example, a variety of modern core banking providers offer modular and flexible systems to replace institutions' antiquated systems and allow them to more easily "plug in" new capabilities such as PETs.

Conclusion

Despite these complexities, the opportunity presented by PETs is large and growing rapidly. Financial institutions today are unable to "see the whole elephant" of their biggest, most pressing shared problems. Privacy enhancing techniques enable institutions to talk – to communicate information about the trunk, the legs, the tails and the ears – without threatening the competitive confidentiality that institutions rely on to retain their edge, or breaching the privacy that customers expect from the guardians of their data.

Appendix

Benefits and limitations of techniques



Technique #1: Differential privacy

Benefits:

- **Allows for manual control over privacy vs. precision:** Adding noise is not binary – more or less noise can be added depending on the institutions' willingness to sacrifice privacy for utility. In the “How it works” example, the surveyor could instead replace the inputs of two or more individuals instead of only one, introducing a greater amount of noise to the calculation. This ensures greater privacy for the 10th individual, but is also less useful to the others in the room in determining if they are overpaid or underpaid, since they are less confident that the number is accurate.
- **It is possible to mathematically measure the privacy leakage:** In reality, the “noise” is added (and accounted for in the summary statistic outputs) through well-defined mathematical formulas and is measured as “differential privacy”. This makes it possible to statistically measure the amount of privacy being leaked by any given output, and institutions can make a conscious choice whether the privacy leakage is acceptable for the value being derived. This level of statistical control makes the technique very customizable to the sensitivity of the data in question.
- **Computationally inexpensive:** Adding noise does not require significant additional computing power for a data-sharing initiative vs. a traditional direct transfer of data. In the “How it works” example, it is not significantly more effort for the intermediary to add a random noise factor to the analysis, then proceed with the calculation of the average as usual.

Limitations:

- **Can be used only on large datasets:** With smaller datasets, it is not possible to add enough noise to protect the privacy of individual contributors while still providing specific-enough information to be useful aggregate statistics. In the “How it works” example, if there were only three colleagues in the room, the surveyor replacing one of the inputs could have a meaningful impact on the calculated average outcome, to the point where it may not be of value at all for individuals seeking to determine if they are overpaid or underpaid.
- **Limits precision:** Adding noise to the inputs, computation of the inputs or the outputs ultimately reduces the precision of the analysis. As a result, differential privacy cannot be applied in specific situations where precise results are critical (e.g. anomaly detection, which can rely on detecting small but statistically significant differences between values).



Technique #2: Federated analysis

Benefits:

- **Minimizes communication costs:** In some cases, especially when large volumes of data are involved, sharing the data itself can become prohibitively expensive. Federated analysis allows for much more concise insights to be shared instead. In the “How it works” example, rather than duplicating the full contents of the reported spam emails from each provider into a central database, only the succinct spam engine insights need to be shared.

Limitations:

- **Requires certain scale of data within each dataset:** Federated analysis assumes that meaningful insights can be derived from isolated sets of data: In some cases, this scale of data may not exist and would lead to limited value being derived from federated analysis. In the “How it works” example, if each individual email provider did not have enough data to independently create useful spam-prevention models, attempting to do so through federated analysis would also not yield any valuable results.
- **Complexity of distributed systems:** Managing a federated ecosystem is significantly more complicated than a traditional centralized database. In the “How it works” example, when institutions are defining their individual spam engines, there are three sets of analysis (one conducted by each company) and no communication between them. When institutions create a centralized database, there are three sets of communication (as the three institutions contribute their data into one database) and one set of analysis (on the centralized database). In the federated ecosystem, there is both the three sets of analysis (as each institution conducts its own analysis) *and* three sets of communication (as they share the insights from their internal analysis). While, as identified in the “Benefits” section, this is low cost in terms of the volume of data being transferred, this communication still introduces additional complexity.



Technique #3: Homomorphic encryption

Benefits:

- **No trust in third parties required to ensure privacy:** Most forms of privacy require some trust in a specific third party (e.g. a certification body). With homomorphic encryption, there is no need for this trust, and data can be shared with a broader set of players. This can allow for more competition and innovation in a market, as providers do not need to go through burdensome certification processes in order to participate in the market and attract customers. In the “How it works” example, when searching for a health analytics provider, Susan does not need to trust that HealthAnalytics Co. will be a good steward of her data or that it has robust security protocols in place; Susan is free to choose them solely based on the quality of their analysis, assuming that homomorphic encryption is in place.

Limitations:

- **Technologically limited:** The technology underpinning homomorphic encryption today is limited either by simplicity or efficiency...
 - Analysis conducted on fully homomorphically encrypted data is orders of magnitude slower than the same analysis on the underlying encrypted data (dependent on the complexity of the calculation). This increased computational cost and latency means that fully homomorphic encryption is applicable only in certain use cases that are not particularly time sensitive. In the “How it works” example, HealthAnalytics Co. would require a significantly greater amount of time to perform any meaningful analysis on the health records provided by Susan than if the data was shared through traditional means.
 - To speed up the analysis, a different type of homomorphic encryption can be used that allows for only one or a few type(s) of operations on the underlying data (e.g. addition or multiplication, but not both). This encryption is known as homomorphic encryption (HE) as opposed to fully homomorphic encryption (FHE). While using HE would significantly increase the speed of the analysis, it limits the depth of insights that can be driven from the data, as different operations cannot be combined. In the “How it works” example, the value of HealthAnalytics Co.’s analysis would be limited to simpler operations, providing less meaningful insights to Susan regarding her health.
- **Verifying results:** Most HE and FHE schemes are not verifiable, meaning that the system cannot provide a proof that the output it has calculated is accurate. As a result, parties using the system must have confidence that the encryption scheme is accurate and has not been interfered with to produce inaccurate results. Verifiable (fully) homomorphic encryption systems are in development, but are even more technologically limited than FHE and HE schemes.



Technique #4: Zero-knowledge proofs

Benefits:

- **Simple to implement:** Zero-knowledge proofs are not mathematically complex and can be integrated into existing systems with relative ease. In the “How it works” example, Peggy and Victor do not need to perform complex mathematics to conduct the exchange of information.
- **Increase security without significant impact to customers’ experiences:** Many other security and privacy measures impede customers’ experiences. For example, two-factor authentication on retail payments would slow down the purchase process (e.g. as customers wait for a text on their phone and then enter it to validate their identity) and increase the complexity of making purchases with a credit card. Zero-knowledge proofs could be integrated into a payment schema without requiring significant additional effort from customers. While in the “How it works” example, the “sample/move-cups/sample-again” process needs to be repeated many times in order for Victor to be able to know with certainty that Peggy knows the difference between the two brands of soda and isn’t just guessing, in reality this interaction would be between computers, happening automatically and much more quickly.

Limitations:

- **Computationally expensive:** While traditional interactive proofs require limited interactions in order for a customer to prove something, zero-knowledge proofs require much more effort. Consider the first proof in the “How it works” example, where Peggy just tells Victor how to differentiate between the two sodas. Compared to the final zero-knowledge proof – where Victor switches the glasses for several rounds, allowing Peggy to sample each drink before and after – the ZKP process requires significantly more effort. While there are non-interactive “zero-knowledge proofs” where Peggy and Victor don’t need to repeatedly communicate, the prover (Peggy) will always need to perform significantly more work in a ZKP system to prove their knowledge.



Technique #5: Secure multiparty computation

Benefits:

- **No trust in third parties required:** Most forms of security and privacy require some trust in a specific third party (e.g. a data analysis entity). SMC removes the need for this trust, allowing for individuals with shared objectives to collaborate with each other. By requiring the consensus of multiple parties, companies do not need to trust the other participants in the data-sharing collaboration. In the “How it works” example, there is no need for any individual to trust other specific participants, but only to trust that the collaboration as a whole is focused on achieving the intended benefits. The threshold amount can be configured to a higher/lower number depending on the level of trust that participants to the data-sharing collaboration have for each other (among other factors).
- **Computationally inexpensive:** Unlike homomorphic encryption, the lack of complex encryption and analysis of encrypted data means that analysis itself can be conducted easily. In the “How it works” example, very little additional effort (in terms of mathematical operations) is required to add SMC vs. the hypothetical case with a trusted intermediary.

Limitations:

- **High cost of communications:** Unlike homomorphic encryption, the cost of communications is significantly higher. In the “How it works” example, conducting the analysis requires each party to go through multiple steps to arrive at a relatively simple outcome; compared to the situation where a trusted intermediary is used, there is much more back-and-forth between participants.
- **High set-up costs:** SMC systems need to be designed and customized to every use case separately. As a result, setting up an SMC system can be expensive and time-consuming. In comparison, differential privacy relies on standardized and generalized mathematical formulas regardless of the data being shared and the analysis being performed.

Further reading

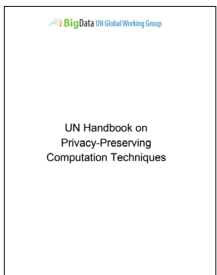
Protecting Privacy in Practice, The Royal Society



Is Privacy Privacy?, Berkman Klein Center



UN Handbook on Privacy Preserving Computation Techniques, Big Data UN Global Working Group



Acknowledgements

Project team:

Jesse McWaters, Financial Innovation Lead, World Economic Forum
jesse.mcwaters@weforum.org

Matthew Blake, Head of Future of Financial and Monetary Systems, World Economic Forum
matthew.blake@weforum.org

Rob Galaski, Global Leader, Banking & Capital Markets, Deloitte Consulting
rgalaski@deloitte.ca

Hemanth Soni, Senior Consultant, Monitor Deloitte, Deloitte Canada
hemasoni@deloitte.ca

Ishani Majumdar, Senior Consultant, Omnia AI, Deloitte Canada
ismajumdar@deloitte.ca

Steering Committee:

Josh Bottomley, Global Head of Digital , HSBC

Pierre-Olivier Bouée, Chief Operating Officer, Credit Suisse

Nick Cafferillo, Chief Data and Technology Officer, S&P Global

Vanessa Colella, Chief Innovation Officer; Head, Citi; Citi Ventures

Juan Colombas, Executive Director and Chief Operating Officer, Lloyds Banking Group

Robert Contri, Global Financial Services Leader, Deloitte

David Craig, Founding Chief Executive Officer and Board Member, Refinitiv

Tony Cyriac, Enterprise Chief Data & Analytics Officer, BMO

Rob Goldstein, Chief Operating Officer & Global Head of BlackRock Solutions, BlackRock

Greg Jensen, Co-Chief Investment Officer, Bridgewater Associates

Prof. Dr Axel P. Lehmann, President Personal & Corporate Banking Switzerland, UBS

Lena Mass-Cresnik, PhD, Chief Data Officer, Moelis & Company

Max Neukirchen, Head of Strategy, JP Morgan Chase

Kush Saxena, Chief Technology Officer, Mastercard

Nicolas de Skowronski, Member of the Executive Board, Head of Advisory Solutions, Julius Baer

Michael Zerbs, Group Head and Chief Technology Officer, Scotiabank

Working Group:

Sami Ahmed, Vice-President & Head of Data, Analytics and AI Transformation, BMO Financial Group

Secil Arslan, Head of R&D and Special Projects, Yapi Kredi

Tim Baker, Global Head of Applied Innovation, Refinitiv

Beth Devin, Managing Director & Head of Innovation Network & Emerging Technology, Citi

Roland Fejfar, Head Technology Business Development & Innovation EMEA/ APAC, Morgan Stanley

Gero Gunkel, Group Head of Artificial Intelligence, Zurich Insurance

James Harborne, Head of Digital Public Policy, HSBC

Milos Krstajic, Head of AI Solutions, Claims, Allianz SE

Wei-Lin Lee, Senior Director, Strategy and Growth, PayPal

Juan Martinez, Global Head of APIs, Identity & Connectivity, SWIFT

Michael O'Rourke, Head of Machine Intelligence and Global Information Services, NASDAQ

Jennifer Peve, Managing Director, Business Development and Fintech Strategy, DTCC

Jim Psota, Co-Founder & Chief Technology Officer, Panjiva (S&P Global)

Nicole Sandler, Innovation Policy Global Lead, Barclays

Annika Schröder, Executive Director, Artificial Intelligence Center of Excellence, UBS

Chadwick Westlake, Executive Vice-President, Enterprise Productivity & Canadian Banking Finance, Scotiabank

Special thanks:

Pavle Avramovic, Associate, Financial Conduct Authority

Jordan Brandt, Co-Founder and Chief Executive Officer,
Inpher

Andrew Burt, Chief Privacy Officer and Legal Engineer,
Immuta

Anton Dimitrov, Senior Software Engineer, Inpher

Mariya Georgieva, Director of Security Innovation, Inpher

Dimitar Jetchev, Co-Founder and Chief Technology
Officer, Inpher

Iraklis Leontiadis, Senior Cryptography and Security
Research Engineer, Inpher

Kevin McCarthy, Vice-President, Inpher

Bhaskar Medhi, Co-Founder and Chief Executive Officer,
Ziroh Labs

Alfred Rossi, Research Scientist, Immuta

Additional thanks:

Derek Baraldi

Mary Emma Barton

Andre Belelieu

Kerry Butts

Alexandra Durbak

Natalya Guseva

Kai Keller

Courtney Kidd Chubb

Abel Lee

Nicole Peerless

Denizhan Uykur

Han Yik

Endnotes

1. <https://www.thomsonreuters.com/en/press-releases/2018/june/thomson-reuters-expands-sentiment-data-to-track-top-100-cryptocurrencies.html> (link as of 9/9/19).
2. <https://www.newswire.ca/news-releases/wealthsimple-announces-partnership-with-mint-577957751.html> (link as of 5/8/19).
3. <https://www.nordea.com/en/press-and-news/news-and-press-releases/press-releases/2019/07-05-14h00-the-collaboration-of-six-nordic-banks-results-in-a-joint-kyc-company.html> (link as of 5/8/19).
4. <https://www.gov.uk/cma-cases/review-of-banking-for-small-and-medium-sized-businesses-smes-in-the-uk> (link as of 5/8/19).
5. <https://www.lenddo.com/products.html#creditscore> (link as of 5/8/19).
6. <https://marketingland.com/survey-58-will-share-personal-data-under-the-right-circumstances-242750> (link as of 5/8/19).
7. Note: We do not intend to suggest that regulators requiring institutions to share data (e.g. PSD2) and regulations seeking to protect customers' privacy (e.g. GDPR) are contradictory. Both forms of regulation are unified in their intent to increase customers' control over how their data is managed and transparency into how it is used. However, they do present a conundrum for financial institutions: They increase the complexity associated with sharing data, and at the same time require institutions to make certain owned datasets (e.g. transaction data) shareable.
8. <https://www.ftc.gov/tips-advice/business-center/privacy-and-security/gramm-leach-bliley-act> (link as of 5/8/19).
9. <https://newsroom.ibm.com/Cybersecurity-and-Privacy-Research> (link as of 5/8/19).
10. <https://www.theguardian.com/news/2018/mar/17/cambridge-analytica-facebook-influence-us-election> (link as of 5/8/19).
11. <http://press.capitalone.com/phoenix.zhtml?c=251626&p=irol-newsArticle&ID=2405043> (link as of 5/8/19).
12. <https://www.blog.google/technology/safety-security/project-strobe> (link as of 5/8/19).
13. <https://techcrunch.com/2019/01/31/aadhaar-data-leak> (link as of 5/8/19).
14. Note: Critically, the use of additional customer data may inadvertently harm individuals if the data is used to enable discriminatory or predatory behaviour. This will be explored in greater detail in the context of artificial intelligence in an upcoming paper by the World Economic Forum.
15. Note: These techniques are not mutually exclusive nor exhaustive. As explored through the use cases, they are interconnected and can be deployed together. Some techniques can be used as intermediary steps to enable other (e.g. zero-knowledge proofs to support secure multiparty computation). There are also other techniques, such as trusted execution environments, that are not included in the scope of this report.
16. https://www.census.gov/newsroom/blogs/random-samplings/2019/02/census_bureau_adopts.html (link as of 5/8/19).
17. Cynthia Dwork, Frank McSherry, Kobbi Nissim and Adam Smith (see note 19).
18. <https://www.microsoft.com/en-us/research/publication/calibrating-noise-to-sensitivity-in-private-data-analysis> (link as of 5/8/19).
19. https://www.apple.com/privacy/docs/Differential_Privacy_Overview.pdf (link as of 5/8/19).
20. <https://github.com/OpenMined/PySyft> (link as of 5/8/19).
21. <https://medium.com/tensorflow/introducing-tensorflow-privacy-learning-with-differential-privacy-for-training-data-b143c5e801b6> (link as of 5/8/19).
22. <https://techscience.org/a/2015092903> (link as of 5/8/19).
23. This concept is explored in greater detail, in the context of health-related data, in a White Paper by the World Economic Forum titled Federated Data Systems: Balancing Innovation and Trust in the Use of Sensitive Data.
24. <https://www.groundai.com/project/applied-federated-learning-improving-google-keyboard-query-suggestions> (link as of 5/8/19).
25. <https://medium.com/tensorflow/introducing-tensorflow-federated-a4147aa20041> (link as of 5/8/19).

26. <https://thenextweb.com/security/2017/12/05/personal-info-31-million-people-leaked-popular-virtual-keyboard-ai-type> (link as of 5/8/19).
27. Note: Federated analysis alone would not necessarily prevent sensitive personal information from being transferred up into the centralized master model. Additional privacy practices (e.g. filters) would be needed.
28. <https://venturebeat.com/2019/06/03/how-federated-learning-could-shape-the-future-of-ai-in-a-privacy-obsessed-world> (link as of 5/8/19).
29. <https://www.techrepublic.com/article/is-homomorphic-encryption-ready-to-deliver-confidential-cloud-computing-to-enterprises> (link as of 5/8/19).
30. <https://crypto.stanford.edu/craig> (link as of 5/8/19).
31. <https://www.theguardian.com/news/2018/mar/17/cambridge-analytica-facebook-influence-us-election> (link as of 5/8/19).
32. <https://dl.acm.org/citation.cfm?id=22178> (link as of 5/8/19).
33. <https://www.wthr.com/article/kohls-cashier-accused-stealing-customers-credit-card-numbers-and-using-them> (link as of 5/8/19).
34. <https://www.itwire.com/security/85702-warning,-card-fraud-could-cost-retailers-us130-billion-over-next-5-years.html> (link as of 5/8/19).
35. <https://www.americanbanker.com/news/how-zcash-tries-to-balance-privacy-transparency-in-blockchain> (link as of 5/8/19).
36. <https://coinjournal.net/blockchain-security-platform-nucypher-raises-us4-3m-cryptofunds-vcs> (link as of 5/8/19).
37. <https://www.unlock-bc.com/news/2019-04-30/nuggets-selected-for-fca-innovation-sandbox-for-blockchain-solution> (link as of 5/8/19).
38. <https://www.ingwb.com/themes/distributed-ledger-technology-articles/ing-launches-major-addition-to-blockchain-technology> (link as of 5/8/19).
39. Two-party computation is a special case of SMC where only two entities are involved in the data-sharing (as opposed to multiparty computation, where a number of entities can be involved). For the purposes of this paper, we will focus on multiparty computation.
40. <https://apps.dtic.mil/dtic/tr/fulltext/u2/a066331.pdf> (link as of 5/8/19).
41. <https://ercim-news.ercim.eu/en73/special/trading-sugar-beet-quotas-secure-multiparty-computation-in-practice> (link as of 5/8/19).
42. <https://www.space.com/5542-satellite-destroyed-space-collision.html> (link as of 5/8/19).
43. <https://www.fbi.gov/stats-services/publications/insurance-fraud> (link as of 5/8/19).
44. Federal Deposit Insurance Company, <https://research.fdic.gov/bankfind/index.html> (link as of 5/8/19).
45. <https://www.bbva.com/en/bbva-launches-first-baas-platform-in-the-u-s> (link as of 5/8/19).
46. Note: Secure multiparty computation would require a more complex and time-consuming set-up than an API-based system, and would need to be customized for each partner. As a result, it is probably feasible only for valuable use cases that include highly sensitive data.

One example of this outside of retail banking would be for hedge funds seeking to evaluate a third-party data provider. The hedge fund would want to know that the data it is looking to purchase would improve the quality of its models, but the data provider would not want to share its proprietary datasets prior to receiving full payment. Using a small sample set or historical data would not accurately represent the real-world performance of the data in the hedge fund's models. A secure multiparty computation system would allow the two parties to compute the impact of the private dataset on the hedge fund's models without exposing the data itself.

47. <https://news.gallup.com/poll/162872/one-three-americans-prepare-detailed-household-budget.aspx> (link as of 5/8/19).
48. <https://benefittrends.metlife.com/us-perspectives/work-redefined-a-new-age-of-benefits>(link as of 5/8/19).
49. <https://www.fraud-magazine.com/article.aspx?id=4294990598> (link as of 5/8/19).
50. <https://www.weforum.org/reports/disruptive-innovation-in-financial-services-a-blueprint-for-digital> (link as of 5/8/19).
51. <https://arxiv.org/abs/1111.5228> (link as of 5/8/19).



COMMITTED TO
IMPROVING THE STATE
OF THE WORLD

The World Economic Forum, committed to improving the state of the world, is the International Organization for Public-Private Cooperation.

The Forum engages the foremost political, business and other leaders of society to shape global, regional and industry agendas.

World Economic Forum
91–93 route de la Capite
CH-1223 Cologny/Geneva
Switzerland

Tel.: +41 (0) 22 869 1212
Fax: +41 (0) 22 786 2744

contact@weforum.org
www.weforum.org



OPINION

Privacy Does Not Pause in Pandemics

BY **SUNNY SEON KANG**

April 6, 2020

When crisis strikes, privacy is too often brushed aside as a competing interest that detracts focus from the greater problems ahead.

But conceding privacy as the first sacrificial right in an emergency means we excuse policymakers from engaging in a careful assessment of the necessity, proportionality and invasiveness of measures that carry long-term consequences.

We must reject the false dichotomy of “lives over privacy” and examine how public authorities can combat the COVID-19 pandemic within justifiable bounds of civil liberties.

Without demands for accountability, we could be turning a blind eye to potential misuses of sensitive data and the accumulation of unchecked information monopolies. Now is the time to keep a vigilant watch over how governments and companies are collecting and processing data — and demand legal and technological reform against systemized surveillance.

It should be very concerning to Americans that the United States does not have a federal privacy law that establishes baseline data ethics during this pandemic.

Without a common denominator of privacy and security standards, risky and short-sighted data practices often fall through the cracks of regulation.

Big Tech is the greatest beneficiary of this patchwork

Big Tech is the greatest beneficiary of this patchwork system, which fosters an overindulgence in user data. And they stand to gain even more if unfettered data collection is condoned through this crisis.

Facebook and Google, in particular, have a long history of privacy violations prosecuted by the Federal Trade Commission. Virtually all of their transgressions involved blindsiding consumers into an uncharted collection and use of their personal data. Now, they want to assist public health policymakers by analyzing this aggregate data for COVID-19 tracing and tracking.

Location data is undoubtedly valuable to epidemiology. Singapore and Hong Kong have both used smartphone data for COVID-19 contact-tracing with measurable success. Still, we must pause to consider the sensitivity of the recycled social media data offered by Facebook and Google, as well as the lack of accountability awaiting these companies if things go wrong.

A dataset as idiosyncratic and detail-rich as someone's real-time location history is highly unlikely to stay anonymized without the assistance of advanced cryptographic privacy safeguards. Research published last year from Imperial College London exposed the inadequacies of mainstream anonymization techniques by accurately re-identifying 99.98 percent of Americans in a dataset scrubbed of personal identifiers.

Worse yet, it may not even take a team of researchers to undermine anonymity. In South Korea, health authorities have inadvertently stigmatized and exposed COVID-19 positive individuals in mass "public safety alerts" that published their GPS movements and demographic information. Replicating these efforts in the U.S. would not only be disturbingly intrusive, but

is likely to be ineffective in isolating the disease at this

is likely to be ineffective in reducing the spread at this stage of a nationwide contagion.

Now is not the time to place blind trust in Big Tech. It is alarming that Facebook and Google are currently bound by FTC consent orders, yet the agency tasked with monitoring the companies' privacy programs have not clarified how a potential engagement with the U.S. government on COVID-19 measures will impact their compliance obligations.

This pandemic brings into focus the failures of the U.S. sectoral privacy system. Mobile location data is not "Protected Health Information" under the federal health privacy law, HIPAA. Since the user would be opting-in to log the location data on their own initiative, HIPAA would not apply even if this aggregate data is directly used for public health analysis.

The tech industry's interference with fundamental rights also poses a constitutional issue. Governmental coordination with Big Tech on national disease control would give these companies the power to set public policy by proxy. This means that important public health decisions could be made inside the black box of private surveillance technologies — operating on potentially biased, unrepresentative, or inaccurate datasets.

COVID-19 presents a compelling interest for the government to act now. Yet, the urgency to accelerate a public health response must be balanced with proper safeguards for civil liberties and privacy.

In the absence of sufficient legal safeguards against data misuse, technical safeguards should be employed to automate data minimization, limited retention, and purpose limitation. Privacy-enhancing technologies that enable distributed computing — such as homomorphic encryption and secure multi-party

computation — can help strike the balance of data

computation can help strike the balance of data utility and privacy whilst curbing data overreach. Built-in privacy is a failsafe where policy and regulations fall short.

People deserve concrete assurances from the government that new data-driven measures tackling the epidemic will not continue to monitor them years into the future.

We are not selfish for demanding answers for privacy right now.

Sunny Seon Kang is Senior Privacy Counsel and Head of Policy at Inpher, a privacy-preserving machine learning company, and advocates for regulatory agencies and private companies to address consumer rights and civil liberties at the intersection of law and technology.

Morning Consult welcomes op-ed submissions on policy, politics and business strategy in our coverage areas. Updated submission guidelines can be found here.

MC/TECH: SUBSCRIBE

Get the latest news, data and insights on key trends affecting tech and tech policy.

SIGN UP

TECH