**JIPS** informing
solutions to internal
displacement

# Can Advanced Data Science Methods Be **A Game-Changer For Data Sharing?**

Innovation Project | Phase 1

2019-2020

# Acknowledgements

The Innovation Project builds on JIPS' exploratory work in collaboration with UNHCR on the use of the statistical disclosure control (SDC) method, which started in 2018. Phase I of the project, focused on the review of data science methods for data sharing which triggered discussions of their applicability in the humanitarian and development fields. Being convinced that these methods have great potential to change the fundamental rules of trust and collaboration between stakeholders, JIPS explored partnerships with like-minded organisations and potential users in the field. The present report is the result of those efforts and discussions.

## Who is JIPS?

JIPS – the Joint Internal Displacement Profiling Service – is an interagency field support service dedicated to bringing governments, displaced persons, host communities and national and international actors together to collaborate towards durable solutions. As a globally recognised neutral broker, JIPS supports collaborative and responsible approaches to data collection as well as its use in internal displacement contexts. Our focus lies in developing national capacities to respond to internal displacement through informed and evidenced-base decisions.

# Table of Contents

# Executive Summary

This document reviews data-science methods for data sharing and discusses their application in the humanitarian and development fields. We believe that these methods have **great potential to change the fundamental rules of trust and collaboration between stakeholders**. To make more efficient use of existing data and gain better insights, we need to look beyond data anonymisation alone and think of a holistic workflow for data sharing and querying.

We argue that by making data and insights safe and secure to share between stakeholders, it will allow for a more efficient use of available data, reduce the resources needed to collect new data, strengthen collaboration and foster a culture of trust in the evidence-informed protection of people in displacement and crises.

The paper first defines the problem and outlines the processes through which data is currently shared among the humanitarian community. Against this backdrop, the report delves into the background for JIPS' engagement in this area of work made possible through the UNHCR Innovation Fund.

It then provides a brief introduction into privacy-preserving methods and their central importance in ensuring safe data sharing, providing an overview of relevant state-of-the-art approaches such as Differential Privacy Algorithms, as well as existing technical solutions for privacy-protected sharing and querying of sensitive, individual-level micro-data. Based on this discussion, we then outline recommendations and next steps for developing a new humanitarian data-sharing framework.

We suggest testing and prototyping the safe extraction of insights from individual-level data, without the need to publish or release the datasets as such (with a prior cost-, resource- and labour-intensive anonymisation). The approach is, rather, to grant protected access to querying the data (which remains on the organisation's servers) through an open algorithm – and, where data is decentralised, through a federated learning-based API – in order to "bring the question to the data".

# 1. Introduction

## Problem statement

Much has changed in the humanitarian data landscape in the last decade and not primarily with the arrival of big data and artificial intelligence. Mostly, the changes are due to increased capacity and resources to collect more data quicker, leading to the professionalisation of information management as a domain of work.

Larger amounts of data are becoming available in a more predictable way. We believe that as the field has progressed in filling critical data gaps, **the problem is not the availability of data, but the curation and sharing of that data between actors as well as the use of that data to its full potential**.

In humanitarian operations and particularly in emergency situations in contexts of forced and protracted displacement, timeliness and access to sensitive individual-level microdata is necessary for planning and carrying out the provision of assistance and the protection of people in distress. The response is informed by data deriving from registration, the profiling of IDP situations, the assessment and prediction of needs and other data on the ground.

Despite the consensus within the community to increase collaboration, data sharing and interoperability of data, the reality remains difficult. The sharing of anything other than highly aggregated results continues to be the exception rather than the rule; where microdata is shared, it is mostly bilateral, ad-hoc, and in many cases informal (e.g. via unprotected excel files), which places sensitive information at a higher risk of disclosure. The discourse on data sharing is also seeing a new trend: an increased use of legal data-sharing agreements that often are set up to comply with normative frameworks such as the GDPR but are rarely directly applicable to the operational reality of data actors on the ground. Therefore, data that can be used to identify needs, plan responses and monitor progress often remains locked away, stored on organisations' servers with a quickly expiring shelf life.

There are a number of barriers to sharing data, but **a prominent one is a lack of knowledge about and a capacity for data anonymisation methods, which are a first step in addressing privacy concerns**. While national statistics offices (NSOs) have built knowledge of handling sensitive census and survey data for years, this is an area of work which has only recently become more relevant for the humanitarian community.

In addition, with growing resources and the acquisition of new skills, the economy of data has also reached the humanitarian domain. Therefore, **data that is kept private to an organisation has become a competitive asset both in terms of business models and with regards to donor-funding and its operational influence in the field**.

The community has come together in the past years, particularly as mandated by the Grand Bargain, to discuss issues of data governance, design pathways for data sharing, promote joint approaches,

expand collaboration, and emphasise the interoperability of data and analysis. If, however, commitments to trust are not underpinned with rigorous use of methods and technology, they will continue falling short of exploring the full potential of humanitarian and development data.

**To make more efficient use of existing data and gain better insights as well as avoid a duplication of efforts, we need to look beyond data anonymisation alone and think of a holistic workflow for data sharing and querying**. With the growing granularity of the data available for analysis to the community, as well as the increased integration of mobility data (such as call record details or financial record data among others), anonymisation alone is not powerful enough to prevent the re-identification of individuals. A model developed by researchers from two European universities suggests that complex datasets of personal information cannot be protected against re-identification by current methods of "anonymising" data, such as releasing samples (subsets) of the information[1, 2].

Yet, anonymisation methods are currently among the most used procedures for preserving privacy, but better results in ensuring privacy and facilitating sharing more effectively and on a larger scale may be gained through the safe Q&A of data. This will require continuous and improved rigour in data collection and storage as well as building the capacity of information management officers and data scientists to apply those methods. We have an ethical commitment to understand and make use of these methods: there is an ethical imperative to not only collect the essential amount of data but also share and use it in the safest and most secure ways possible.

Therefore, we believe that it is critical for the humanitarian community to build its knowledge of and capacity for using new, more advanced methods. In our context, we can identify three problems in data sharing that specifically concern microdata[3]:

1. **Microdata holds highly sensitive information** that can put people at risk. Privacy and data protection are therefore paramount but also difficult.
2. The collection of this type of data is a **resource-intensive, time-consuming endeavour** that can put a strain on respondents.
3. A large amount of data collection is undertaken, but **much of it is not shared and analysed**. It simply sits there, stored away across organisations' servers, largely inaccessible for analysis to the larger community.

## The JIPS innovation project

This review is the output of the first phase of the JIPS Innovation Project which was made possible through a grant from the UNHCR Innovation Fund. The review focuses on the exploration of innovative solutions to incentivise microlevel data sharing.

---

[1] Lomas, N. (2019, July). Researchers spotlight the lie of 'anonymous' data. TechCrunch. Retrieved January 29, 2020, from http://social.techcrunch.com/2019/07/24/researchers-spotlight-the-lie-of-anonymous-data/
[2] Rocher, L., Hendrickx J.M., and de Montjoye Y.A. (2019). Estimating the success of re-identifications in incomplete datasets using generative models. *Nature Communications*. 10(1), 1–9, https://doi.org/10.1038/s41467-019-10933-3
[3] The report includes a glossary Terms & Definitions where technical terminology is explained

The innovation project builds upon JIPS' exploratory work in collaboration with UNHCR on the use of the statistical disclosure control method. Within the scope of that work, 45 National Statistics Offices were contacted through the Expert Group on Refugee and IDP Statistics (EGRIS) and asked about their practices of anonymisation and release of microdata. Among those who responded, the following key take-aways emerged from the study:

> **Microdata** is released as public or as scientific use files (licensed files), through data centres or through secure remote access.

> **Public use files (PUFs)** are "created to allow the general public to get familiar with statistical microdata files. The files are prepared in such a way that individual entities cannot be identified. This goes with a loss in information value", as expressed in one response to our study.

> **Scientific or licensed use files** are microdata files suitable for research, as more of the original information has been retained. The disclosure risk is therefore higher than would be tolerated in a PUF, and researchers will typically have to go through an application procedure and sign a confidentiality agreement in order to get access to physical copies of the datasets.

> **Data centres** are venues where researchers and others with the relevant accreditation can access selected microdata files of a national statistical office. The visitors will typically have to use the offline computer facilities that are provided at the venue, and they cannot download, copy or in any other way take microdata with them when they leave. Only output from analyses can be taken outside, but it will often be subject to a check of confidentiality before the user is allowed to do so. Data centres are usually run by NSO staff. Some NSOs also offer online platforms for secure remote access, where users can send queries to non-anonymised databases and view the output, with or without being able to view the actual data.

Some NSOs offer all the above options, whereas others rely on only one or a few. One NSO replied that microdata is currently not being disseminated at all, but that the office is exploring tools for anonymisation, aiming to disseminate microdata in the future. Common for all NSOs reviewed in this project is that their dissemination approaches and their practices for managing the risk of disclosure have been shaped by national legislation, which dictates what can be shared, with whom and in what form.

Building on the above precursory research, JIPS explored partnerships such as OPAL/Data-Pop Alliance, Johns Hopkins University's Applied Physics Lab, Flowminder, the UN OCHA Centre for Humanitarian Data and a number of potential users in the field, including NGOs and Governments in the Northern Triangle of Central America, as well as the Government of Colombia's Victims' Unit.

# 2. Disclosure vs. Analytical Value

In the past, microdata was mainly shared by researchers, but it has become increasingly common also in the humanitarian sector. Regardless of the agents involved, whenever microdata is shared, a sensitive trade-off needs to be balanced: **sensitive individual-level microdata needs to be analysed without putting people at risk**. Therefore, the possibility for re-identifying respondents and the risk of disclosing their identity should remain low, but the analytical value of the data should be maintained. The following section gives a brief overview of additional methods seeking to address this issue.

## Introduction to privacy

Privacy in our context touches upon **two main concepts: privacy of the individual or group, and data protection or security**. In different terms, we are interested in (1) protecting sensitive, personal information that could reveal the identity or identifying attributes of individuals or persons at risk and (2) protecting or securing this data at each step of the way against processing, query and analysis through unauthorised access.

We need to first briefly explain the concepts of **personal information and sensitive information**. Personal data is generally understood as anything that can identify an individual. Here we can distinguish, on the one hand, direct or explicit identifiers, which refer to unique attributes that clearly identify an individual, such as name or driver's license number. These reveal the identity of a person and must be anonymised or concealed before use in order to ensure privacy.

On the other hand, we can have **indirect or quasi-identifiers**, such as date of birth, address, job title or postcode. These may lead to the identification of an individual when combined with other information or used in a specific context, and they also need to be anonymised or concealed before use. In addition, certain information is considered sensitive. Sensitive attributes include, for example, salary, religion, positive or negative test results, health status and preferences that can be used to pressure an individual, such as political affiliation or sexual orientation. Simply redacting direct identifiable attributes such as a person's official name or ID does not work because they can be identified through other attributes or characteristics that can be revealed with other means, such as by cross checking their profession, age and education in another dataset such as granular geographic location or mobility profile. With an increasing number of sources collecting personal data, this is fairly easy and affordable to do.[4]

---

[4] Narayanan, A., and Shmatikov, V. (2008). Robust De-anonymization of Large Sparse Datasets. 2008 IEEE Symposium on Security and Privacy, 111–125, https://ieeexplore.ieee.org/document/4531148

# Methods for data anonymisation

Anonymisation refers to hiding or modifying the identifiers of an individual in a dataset so that the sensitive information regarding that individual cannot be tracked back to them. There are various techniques for anonymisation from simple steps such as removing or redacting directly identifiable attributes (name, email, IDs), to more advanced data-clustering methods such as k-anonymity, l-diversity and t-closeness (see below).

**Pseudonymisation**

It is distinct from anonymisation and refers to the processing of personal data in such a manner that the data can no longer be attributed to a specific individual without the use of additional information, provided that such additional information is stored separately.

**K-anonymity**

After the removal of direct identifiers in a dataset, the quasi-identifiers are grouped in such a way that an individual cannot be identified without the dataset losing the overall essence of the information (see example below). The technical requirement of this approach is that a list of data is partitioned into a certain number of classes (also called equivalence classes) in such a way that for each row of information, there are k-1 other rows that have the same value for the quasi-identifier attributes. This approach protects against identity disclosure but does not prevent the disclosure of sensitive attributes as it is subject to homogeneous and background attacks.

*Example: Suppose this is the original table:*

| Name | Age | Gender | State | Disease |
|------|-----|--------|-------|---------|
| Person X | 30 | M | State C | Cancer |
| Person Z | 70 | F | State A | TB |
| Person Y | 25 | M | State B | TB |
| Person A | 50 | F | State B | Heart Attack |

The following table is the modified table with *2-anonymity* with respect to attributes *Age* and *Gender*. We can see that for each pair of age and gender, the same age and gender values are present in at least two rows. Here, an individual with certain age and gender cannot be identified, but the original information is not lost either.

| Name | Age | Gender | State | Disease |
|------|-----|--------|-------|---------|
| * | 20-30 | M | State C | Cancer |
| * | 50-80 | F | State A | TB |
| * | 20-30 | M | State B | TB |
| * | 50-80 | F | State B | Heart Attack |

## ℓ-diversity

This approach aims to solve the issues of k-anonymity pertaining to its susceptibility to homogeneous and background attacks. ℓ-diversity guarantees that for each quasi-identifying group (e.g., "artists" for painter, dancer, etc., and "age range 30–40" for values such as 30, 35, and 32), there are at least ℓ distinct values for sensitive attributes. Mathematically, it requires that the distribution of a sensitive attribute in each equivalence class has at least ℓ "well represented" values. However, this does not prevent attacks that are based on knowledge about the global distribution of these sensitive attributes. And, in some cases, it offers more information gain to attackers because of the nature of distribution of the sensitive attribute values.

**t-closeness**

This approach aims to overcome the vulnerability of ℓ-diversity by ensuring that the distribution of sensitive attributes in any equivalence class is closer to the distribution of the attribute in the overall table. Technically, for a given quasi-identifier group, it is guaranteed that its distribution is bound by t against its corresponding distribution on the entire dataset.

## Introduction to privacy-preserving mechanisms

To address the issue of protecting individuals' privacy while increasing access to and insights generated from collected data, this section introduces different ways in which data is currently shared in the humanitarian context and various approaches to how privacy can be protected, ranging from anonymisation, statistical manipulation or statistical disclosure control methods, to approaches based on machine learning.

There are various privacy-preservation mechanisms that enter the workflow of data management, analysis and sharing. An important concept that contributes to privacy protection in humanitarian data sharing is privacy by design. **Privacy by design refers to an approach to designing and structuring software in which information is not exposed to anyone else, in or out of the system. In such systems, privacy is treated not as an added feature but equally crucial as the system architecture**. The information remains secure from the time it enters the system to the time it is destroyed properly.

Implementing privacy building blocks in an existing system is more difficult and costly than designing a system with inherent privacy. When the latter is done, the workflows are studied thoroughly from start to end, weak points are identified and only then the system is designed and built. Various approaches such as hashing, encryption and statistical manipulations are applied in the parts of the process where security and privacy are of key importance.

The methods of applying data privacy are described in brief below.

### Encryption

Encryption is the process of encoding a message or information in such a way that **only authorised parties can access it**. There are key(s) involved in encryption which are used to encode and decode the information. It is different from hashing in that the original value can be recovered with an appropriate key, while in hashing there is no way to recover original information once hashed.

### Hashing

Hashing is a process of generating a value of fixed length from input data with a given value using a hash function which has a peculiar property: the hash value of original input data can be calculated easily, but given a hash value, its original value is computationally very expensive to calculate. **Hashing is generally used when we want to hide or obfuscate the original data.** For instance, given an arbitrary hashing function (here, the MD5 method), the hashing results

of the strings "Climate Change" and "Zero Hunger" are equal to 70dec0af9355c68c34ec5e27e957b6d8 and f8026977e040a3cbf6a94b673d369c7e, respectively. While encryption is a codification of the original values, in hashing, there is no way to recover the original value. Therefore, even if the hashing function is known to an intruder, retrieving the original strings from the hashing results is computationally extremely expensive.

### Statistical manipulations

The **original information is transformed based on the general statistical distribution of the data**. Consider [(5, 10), (3, 16)], which can be replaced by [(3, 5), (8, 15), (1, 19)] so that the original data is hidden but has the same statistical properties. In such approaches, it is not required that the original values be extracted; we are interested only in the overall statistical property of the original data.

### Statistical disclosure control

Statistical disclosure control (SDC) techniques can be defined as a set of methods for **reducing the risk of disclosing information on individuals, respondents, businesses or other organisations**. Such methods are only related to the dissemination of information and are usually based on restricting the amount of or modifying the data released. There are different ways that data can be treated in SDC. The choice between non-perturbative and perturbative anonymisation methods in SDC can limit the types of analysis that can be performed on the anonymised version of the dataset without additional information of how the original data has been altered.

## Machine-learning approaches in data sharing

When it comes to machine learning, there are even more factors to consider than simply preventing identity disclosure. Ensuring the safe handling of data as it is ingested into and used for either machine-learning models or in complex workflows and shared, is key. In the following, we review common approaches in the machine-learning literature that aim to ensure the preservation of data privacy when a machine-learning algorithm is trained on and applied to the data.

### Homomorphic encryption (HE)

Homomorphic encryption approaches privacy by **encrypting the source data even before it is processed**. This approach, however, differs from normal encryption in the sense that the original result can be restored by applying decryption. The function used to encrypt the data is called homomorphic function, and it preserves the information content of the original data. This method is regarded as the most secure privacy preserving mechanism of all.

As part of the Innovation Project, JIPS had the opportunity to discuss differential privacy with experts of the research centre Applied Physics Lab, a research centre affiliated with the Johns

Hopkins University. One of the suitable ways discussed for enabling safe and secure microdata querying and sharing was state-of-the-art Homomorphic Encryption Computing Techniques with Overhead Reduction (HECTOR). While it is a relatively new field of exploration even in academia, this approach would combine advanced cryptographic techniques with secure multiparty computation allowing for a simple and easy use of cryptographic techniques to perform secure distributed computations on de-centralised data sources. In phase II of the Innovation Project, we will further explore with experts from APL how these methods could be made applicable to our domain.

### Secure multiparty computation (MPC)

This approach is suited for cases where multiple data providers are not willing to share their data but want to collaborate in building a common statistical or machine-learning model. In the MPC approach, **computations are performed on encrypted inputs from different parties**. No party has any knowledge about the actual data of the others, and they only access the output of the machine-learning model. There are several protocols for MPC based on secret sharing and garbled circuits. In secret sharing, a party splits its data and provides it for others, never revealing full input values. In garbled circuits, the computation is expressed as a Boolean circuit which consists of basic gates such as AND, OR, NOT. MPC also ensures correctness and provides high efficiency.

### Order-preserving encryption (OPE)

This is primarily used in **sharing encrypted databases when specific data selection and filtering operations can be applied**. In contrast to the previous approaches, in OPE the order of the encrypted data is preserved, that is, if the data points $a$ and $b$ have a certain order ($a$ should appear before $b$), encrypted forms of $a$ and $b$ also preserve this order. This allows for the application of specific operations of the encrypted data (such as MIN, MAX, >, =, <), without access to the original data.

### Differential privacy (DP)

Instead of employing data encryption as in the previous methods, this approach **preserves data privacy by adding statistical noise to both individual and aggregated data**. Depending on the data and the statistical or machine-learning aggregation method, the added noise can be defined in various ways. For instance, the noise can be a random number based on a predefined probability distribution such as Laplacian. Such manipulated data can then be shared with a third party to implement and learn machine-learning models. While the third party does not have access to the original data, it is aware of the probability distribution, enabling the required processing to be performed. There are two methods for DP: local and global. In *local DP*, noise is defined for and added to each individual data point, while in the *global DP* approach, the same noise is added to all data points. DP is highly efficient, and when applied correctly, it provides a precise method of privacy preservation.

## Authentication methods

Authentication is a way of identifying the user. Users are authenticated in a system when they prove their identity in that system. There are basically three methods of authentication, which are briefly explained below.

### HTTP basic authentication

In this simplest method of authentication, users send their encoded username and password in an HTTP header for each request. It is very simple because there is no need for cookies, session IDs or other stateful information as the user sends their credentials in every request. However, this method is susceptible to man-in-the-middle attacks and, if not used with SSL (Secure Socket Layer), the credentials are completely exposed.

### API keys/tokens

Somewhat simpler and more secure than basic authentication, this method involves a unique generated value that is assigned to each first-time user, signifying that the user is known. When the user attempts to re-enter the system, their unique key is used to prove that they are the same user as before. The use of keys makes this method very fast. To make it more secure, each token can have an age of validity.

### OAuth

This is the most secure method. In OAuth, when users want to log in to a system, the system requests for authentication, usually in the form of a token. The request is forwarded to an authentication server which will either accept or reject the request. If accepted, the user is provided with a token which is then presented to the system. The validity of the token can then be checked at any time independently of the user, and it can be used over time with strictly limited scope and age of validity.

## Data sharing in the humanitarian context

Recent years have witnessed a shift towards increased humanitarian (micro)data sharing, deeper humanitarian collaboration and a stronger engagement with the private sector. We have also seen a heightened general focus on the responsible use of data and the importance of privacy protection within our field and beyond, supported by organisations such as the UN OCHA Centre for Humanitarian Data, UN Global Pulse, the Data-Pop Alliance and a variety of informal communities of practice.

Sharing data among the humanitarian community is currently based on trust agreements, extensive contractual sharing agreements and the use of open-data platforms such as the Humanitarian Data Exchange service. The recent addition of sdcMicro has complemented and supported those data sharing mechanisms.

In addition to advanced privacy preserving methods, there has been further research on the governance of data sharing from actors engaging with the private sector and seeking to extract insights from sources such as Call Detail Records and Mobility Data for Good. The Data-Pop Alliance has particularly built on previous references looking into Data Sharing Paradigms for Good.

## Humanitarian data sharing paradigms

Generally, we can distinguish four different (humanitarian micro) data sharing paradigms that approach privacy protection from different angles, with implications on what kind of analyses they permit, how secure they are, what areas of application they have, and how collaborative they are. They exist largely on a spectrum of privacy protection: they either provide secure access to pseudonymised data that can be queried by users, or they grant direct access to aggregated data and the provided results. They also differ in how collaborative and scalable they are: they may allow sharing from one, few or many entities and datasets to a varied number of users. The sharing paradigms are:

1. **Limited release of data**, which centres primarily around releasing prepared data and aggregated results for public use or data challenges (one-to-many sharing);
2. **Remote access** to (pseudonymised) data, based on contractual agreements and trust (one-to-one/one-to-few);
3. **APIs and open algorithms**, allowing a limited number of question-and-answer operations on datasets through decentralised access (few-to-many);
4. **Precomputed indicators and synthetic data**, which offer insights through information generated for a specific purpose or audience but no actual access to the data itself (one-to-one/one-to-few).

The Data-Pop Alliance introduces a fifth option, **data collaboratives**, in which institutions come together to set up managed data-sharing platforms for collaborative use (few-to-few/many-to-many). However, the latter is not a paradigm as such, as it relies on one of the other four approaches in terms of the technical realisation of the actual method of data sharing; rather, it tackles the governance and collaboration around the use of the data. Nonetheless, it facilitates access and collaboration, and addresses the limitations of the four main sharing paradigms that are resource intensive, burdensome and time consuming. However, most of what is shared in data cooperatives are typically aggregated insights and not the raw or even pseudonymised data.

For a description of this classification of data sharing paradigms, either see the Data Sharing Paradigms Section in the Annex or consult "Sharing is Caring" (2019) by the Data-Pop Alliance, which elaborates on the concept in a discussion on using private, personal data derived from call detail records (CDR) that are collected by telecom operators to inform humanitarian decision making and using (private) data for (public) good initiatives.

For our specific context of wanting to enable the safe access to and use of sensitive individual level microdata for humanitarian data sharing and collaborative analysis, we focus on the third paradigm: APIs and open algorithms. This paradigm is the best suited and most promising starting point for

exploring data sharing in the humanitarian context, as it evades the limitations of the other three while allowing the integration of powerful privacy-protection mechanisms in a scalable way that is suitable for sensitive microlevel data. However, we recognise that this option should ultimately support what the Data-Pop Alliance envisions as "data collaboratives" in the humanitarian and development context.

## Steps of the data sharing workflow

Regardless of the paradigm used or the solution provided, all workflows contain core steps or stages at which privacy and sharing need to be considered. These steps are the lens through which we look into existing solutions and frame what a proposed solution will need to take into account at each point.

### Collection of raw data

Raw data is basically the primary data collected from various surveys, evaluations and assessments or data shared by affected populations with humanitarian actors through e.g. complaints or feedback mechanisms. (Primary data is defined as data collected for a specific purpose; any data collected for a different purpose is referred to as secondary.) Raw data contain individual- or group-level data, which might contain sensitive information that could harm the population.

### Anonymisation

When ingesting or storing the raw data in any kind of system designed for analysis and evaluation of data, privacy considerations are a must. The design of those systems must, therefore, include privacy preservation as an inherent element. The very first and crucial step would be to remove information identifying individuals or groups before the data is even let into the system. The process of hiding or redacting information that could lead to identity theft is called anonymisation or pseudonymisation and is considered the foundation of privacy by design.

### Ingestion and storage

After the data is collected and anonymised or pseudonymised, it needs to be stored within the system for further manipulation and analysis. But even before it is merely stored, some pre-processing is done on the data without information loss so that it is usable within the system.

### Computation (analysis)

The ingested data needs to be operated on different computations in order to yield valuable insights. The computation can consist of aggregation or filtering operations, classification of the data, other types of analysis by an analyst or machine learning (ML) or natural-language-processing (NLP) operations.

**Anonymisation of the performed computation**

The privacy of individual-level data is preserved by another step of anonymisation, applied after or during the computation. This step ensures the preservation of the privacy of low-level data, when high-level knowledge is provided to users.

**Query & response**

The computed data is then made accessible to data users. This can be achieved in ways such as querying the system and exposing the data through APIs.

**Access/authentication**

Upon each query, the level of access to the data, the execution of the computation, and the provided results can be finely controlled and restricted by authenticating users.

## Available solutions for the humanitarian sector

In the humanitarian sector, microdata refers to **data on the characteristics of a population that is gathered through exercises such as household surveys, needs assessments or monitoring activities.** Statistical disclosure control techniques can be defined as a set of methods that aims to reduce the risk of disclosing information on individuals, businesses or other organisations. Such methods are only related to dissemination. They are usually based on restricting the amount of or modifying the released data and are made available through different tools and platforms.

**sdcMicro**



*Figure 1: sdcMicro interface with variables and parameters*

sdcMicro is an open-source R package (https://cran.r-project.org/package=sdcMicro) for assessing, anonymising, and re-evaluating microdata. sdcMicro provides technical specialists with an understanding of the context where they the risk of disclosure in microdata can be managed and controlled.

The package has been created by the International Household Survey Network / World Bank (https://www.ihsn.org/software/disclosure-control-toolbox) and is used by numerous national statistical offices around the world. It is regularly updated to include new anonymisation algorithms (https://www.ihsn.org/sites/default/files/resources/ihsn-working-paper-007-Oct27.pdf). The document also includes a graphical user interface to facilitate initial user training.

*Figure 2: A workflow for applying common SDC methods to microdata*

## FlowKit (Flowminder)

Flowminder is a non-profit organisation with a mission to improve public health and welfare in low- and middle-income countries. Flowminder collaborates with governments, inter-governmental organisations and NGOs and works on **collecting, aggregating, integrating and analysing anonymous mobile operator data, satellite data and household survey data**. The work of Flowminder focuses on the collection and analysis of vast amounts of data that is considered individual-level and private. The analysis provides a mapping of distributions and characteristics of vulnerable populations in low- and middle-income countries.

Flowminder provides these insights and tools for application in other contexts through its FlowKit platform, which analyses call detail records (CDR) collected and provided by mobile network operators (MNOs). Although CDRs are primarily used for billing purposes, the design of FlowKit allows for the extracted data to be used and useful for disaster response, precision epidemiology and transport or mobility. Since CDRs constitute highly sensitive data, FlowKit is designed with state-of-the-art privacy protection in mind – operating within the API/open algorithms paradigm.

There are three barriers to privacy compromisation in the FlowKit:

1. **Ingestion**: Data is pseudonymised prior to it entering the system. Pseudonymisation is done separately before ingestion and not within FlowKit as the system is not intended to have access to the non-pseudonymised data. The data is received whenever available, and its pseudonymisation is the first, crucial step of data privacy. Thus, nobody sees the personal data. It is then translated to the expected structure and saved in the database.

2. **Computation**: The computations consist of aggregation, and results consist of the aggregated data. This is where different algorithms are run on the data collected from CDRs. The algorithms are mostly aggregation algorithms. There is a fine-grained control on what type of computation is performed for a user. For instance, some users can be restricted to view and use the results of a certain computation directly but be allowed to indirectly use the result of another computation that they can only view. Since the computations comprise aggregation, the possible re-identification of an individual remains less likely. However, the re-identification of statistically limited (small) data is still possible. The FlowKit API mitigates this risk by redacting any rows in the aggregated outputs that represent very few observations.

3. **Well-authenticated and finely controlled APIs:** The parties which access FlowKit data and APIs require an authentication which is JWT (JSON Web Token) based. This also serves for tracking which aggregation or computation is delivered for which user. The FlowKit API logs all access requests for auditing purposes. Each token has a time limit and can be checked for access exposure for an indefinite amount of time.

## OPAL (Data-Pop Alliance, MIT, WEF, Vodafone, Imperial College London)

OPAL (from "open algorithms") is a non-profit socio-technological innovation project developed by a group of partners around the MIT Media Lab, Imperial College London, Orange, the World Economic Forum and the Data-Pop Alliance. OPAL aims to unlock the potential of private-sector data for public good and is based on the API/Open Algorithm Paradigm. OPAL operates by "sending code to the data" in a safe, participatory and sustainable manner. OPAL's mission is to apply open algorithms and safe and fair technological and governance systems for better decision making to support sustainable development goals around the globe.

The core of OPAL consists of an open and secure platform and algorithms that can be run on the servers of partner companies behind their firewalls. The following are the key components in OPAL:

1. **Q&A client API**: This is the public interface of the OPAL platform.
2. **Authentication**: This component allows access only to authorised users, who can query and run algorithms on the platform.
3. **Controller/scheduler**: This runs and schedules the computation for OPAL requests.
4. **Data ingestion**: Data from various sources and providers are imported into the OPAL database.
5. **Database**: This is where the imported data and results of computations are stored.
6. **Aggregation service**: This adds noise to the data and implements other redactions to ensure the privacy and security of sensitive information.

7. **Cache**: This component stores the computed values, which can be further used to speed up the computation of popular algorithms.
8. **Auditing service**: This is where every authorised request received by the Client API is logged.

In the question-and-answer approach, an information or knowledge need is devised by a user and sent to the system. The system, after considering the privacy concerns of the information need and the access level of the user, provides an answer accordingly.

Questions are defined by an algorithm in Python, specifying which computation to run on each record and how to aggregate them. OPAL uses the MapReduce model to easily process and aggregate large-scale datasets. The following three steps are taken when executing a question:

1. **Pre-processing**: Data is pre-processed to make it consumable by mappers. Generally, this step includes listing and importing labels.
2. **Mapping**: Each individual entry is mapped to labels.
3. **Reduce**: Aggregation across individuals takes place and a single number or distribution is returned for every label. This step includes various aggregations such as counts, sums, average and standard deviation.

## Data privacy

Disclosure risks are managed using a combination of server-side security, pseudonymisation, fine-grained authorisation for algorithms, limiting private information to be stored or exported and properly safeguarding data. Data privacy is accomplished in the following three steps:

1. **Pseudonymisation and Safe Answers**: All data that contains sensitive information is pseudonymised, and access to answers are limited.
2. **OPAL API**: APIs have fine-grained control over who can make which API calls and computations.
3. **Network security**: Use of secure protocols and design.

## Humanitarian Data Exchange (HDX), the UN OCHA Centre for Humanitarian Data

The Humanitarian Data Exchange (HDX) is an open data-sharing platform managed by the Humanitarian Data Centre of the UN Office for the Coordination of Humanitarian Affairs (OCHA). Launched in July 2014, its goal is to make humanitarian data more available and easier to find and use for analysis. It operates on the premise of the limited-release data-sharing paradigm.

HDX hosts more than 17,000 country-specific or world-wide datasets. More than 250 different organisations share their data through the platform. Data can be shared in three ways: publicly, privately and by request via HDX connect, for which only metadata is provided publicly. HDX does not allow personal data or personally identifiable information (PII) to be shared: all data shared through the platform must be sufficiently aggregated or anonymised to prevent the identification of people or

harm to affected populations and aid workers. According to HDX's policy, microdata should not be shared publicly unless the risk of identity disclosure is below an acceptable level. What an acceptable level of disclosure risk is has currently not been defined – it will have to be specified and modified as more experience is accumulated in the humanitarian community and be established on a context-specific basis.

The UN OCHA Centre for Humanitarian Data uses statistical disclosure control (SDC) processes to mitigate the disclosure risk for datasets shared on HDX. These processes aim to assess and lower the risk of a person or organisation being re-identified from the analysis of microdata.

The SDC process uses the open-source sdcMicro tool and includes an assessment of the disclosure risk for the microdata shared through HDX and a warning that is issued to the contributor in case the disclosure risk is high. The Centre for Humanitarian Data applies the sdcMicro tool to datasets for which the risk of disclosure is above the Centre's maximum threshold.

While HDX serves the broader community, organisation-specific and internal sharing platforms are also in use, such as the UNHCR Raw Internal Data Library (RIDL), which has been operational since 2018 with the aim to create a globally-supported, centralised, and secure data repository that ensures that UNHCR is able to use collected raw data to its full potential, that the data can be preserved for future analysis and use, and that it is available externally to operational partners, project stakeholders or academia. It also relies on the sdcMicro tool for data anonymisation.

# 3. Exploring a Solution

A data-sharing system in the humanitarian context should adhere to the highest standards of privacy protection for sensitive (individual, household, group-level) microdata without the need for time-consuming anonymisation, extraction and transmission from data owners to those interested in accessing insights from it. The solution should allow for a collaborative approach and enable safe few-to-few or multiple sharing.

Privacy protection needs to be guaranteed reliably for insights generated from personally identifiable information by using open algorithms for decentralised, secure access to performing queries and receiving results through an API service.

The question-and-answer (Q&A) paradigm to data sharing provides a possible solution to these needs. The Q&A paradigm entails few-to-many access, which allows the use of private data without the data ever leaving the premises of its owners. In this paradigm, the owners provide APIs through which the users can either access pre-computed indicators or richer answers prepared by the owners. This paradigm has also been used in the GSMA Big Data for Social Good Initiative, the FlowKit system, and the OPAL project. In the Q&A paradigm, the data remains securely with the data owner, addressing

the concern of data being misused "in the wild". Moreover, users get more flexibility as they can exploit the knowledge provided in the data, by asking curated questions through APIs. The APIs allow data owners to enforce any level of authentication on access to the data.

**As part of phase II of the JIPS innovation project, we will work with the Flowminder team using components of FlowKit within a prototype of an advanced humanitarian data-sharing solution**. The FlowKit architecture provides an open-source platform for data ingestion and anonymisation, data aggregation, and Q&A APIs, all under the umbrella of a granular authentication system. Given that FlowKit is domain-specific and has a single-server design, we will consider a system based on Substra instead, with adapted core elements of FlowKit, if the final selected use-case for the prototype is based on a de-centralised source. The addition of a FlowKit-style API that permits controlled Q&A access to users outside the network would add useful functionality to the system. The remaining question is the nature of the algorithms that will be used within the system, and how FlowKit's API/permissions can be adapted to provide access to these algorithms or insights.

In addition to the components mentioned so far, which are already available in FlowKit's architecture, an important extra component comprises privacy-preserving methods for machine learning and aggregation. This component secures the preservation of data privacy, even if the extracted knowledge from data is provided in aggregated statistical form. The implementation of this component will expand the capabilities of the FlowKit system itself and contribute to the open-source community in the humanitarian sector. Any fully developed and quality-controlled prototype as part of this innovation project will be open source and made available to the humanitarian community.

The FlowKit system and the discussed solution is based on a centralised approach, where all data is available on one server. In the case of a decentralised data source, the Substra platform, provided by the Substra Foundation could be an effective complement to the FlowKit system.

Substra is a distributed system built upon data locality, traceability and decentralised trust. In this system, data remains on the premises of each organisational node, and the result is achieved through computations on each node and a final aggregation. To enable this, Substra's framework uses a secure multiparty computation network to execute the secure computation of models. This component consists of nodes and interconnections, where the nodes are the individual organisations that own the private data. Each node is willing to build an individual model on the private data but is unwilling to share the data itself. Each node is connected to the others via the network, through which only computations and models flow. In this framework, the distributed nodes tracking the computations provide computational transparency and traceability, which also flow through the network and to which each node agrees upon. Besides the multiparty computation network component, communication between the network and the "outside world" is enabled by APIs and the authentication layer. A scheme of the proposed decentralised system is shown in *Figure 3*.
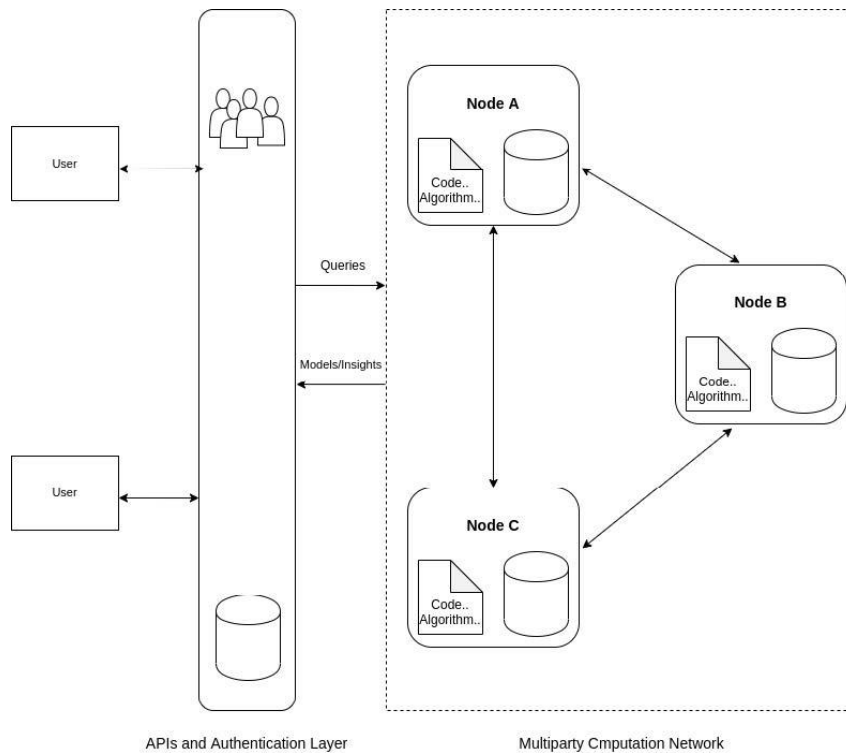
*Figure 3: Proposed system architecture*

A decentralised, or distributed, approach to data sharing through the mechanism described above would allow the safe querying of multiple data sources such as demographic and socioeconomic data located at national statistics offices or with development actors as well as data on current or protracted needs of the affected population.

# 4. Conclusion and Next Steps

We set out to understand what could contribute to changing the game of how humanitarian data is "shared". To tackle the humanitarian data-sharing problem, we have explored current state-of-the-art technical solutions that can be leveraged to unlock the largely untapped potential of existing microdata, in order to provide better assistance for people at risk while still ensuring that both data and its privacy are protected and that the data that has been collected is also used.

The approach we are proposing is to build capacity and test the safe extraction of higher-level insights from individual-level or group-level data without the need to publish or release the datasets as such, while addressing the anonymisation of the data, learning from the experience that already exists in using call detail records.

What we propose instead of publishing the data as such is granting protected access to querying the data (in a decentralised server) through an open-algorithm API based on federated learning in order to "bring the question to the data" and provide secure, privacy-protected individual-level insights back to authorised users, thus skipping the step of moving the dataset or bringing sensitive data to the questioner.

In the next phase, we aim to test the recommended technical solutions and in collaboration with partners develop a prototype model for the safe sharing of humanitarian individual-level microdata. The implementation of the next phase will be realised by JIPS and UNHCR in close collaboration with a range of strategic partners, including Flowminder, the Johns Hopkins Applied Physics Lab, and the Humanitarian Data Centre.

# ANNEX:

## JIPS Innovation Project Timeline



## Terms & Definitions

**Anonymisation –** Techniques aiming to ensure that datasets containing personal data are fully and irreversibly anonymous so that they do not relate to an identified or identifiable natural person, or that the data subject is not (any longer) identifiable.

**Data –** Facts and statistics collected for reference or analysis.

**Data** (various states) **–** One can distinguish four different states of data in this context.

> **Raw data:** Data that has not been modified in any way after being collected. Although containing a lot of information, both relevant and irrelevant, it is not yet ready for consumption to gain insights.

> **Anonymised or pseudonymised data:** Raw data is prone to identity disclosure and privacy breaching as it contains direct identifiers and quasi-identifiers. Different anonymisation or pseudonymisation techniques are applied to prevent the identification of individuals or groups.

> **Processed data:** After the raw data has been collected, it is then processed and made suitable for further use cases. This includes removing irrelevant information, transforming the raw data, validating the data and removing or complementing incomplete information.

> **Aggregated or computed data:** Processed data is simply cleaned data that is made fit for other uses and thus provides no valuable insights. Different computations are performed on the processed data to make results and insights available to end users.

**Data analytics** – The practice of examining data through qualitative and quantitative analysis and research with aims such as gaining insights, identifying behavioural patterns, drawing conclusions, or improving decision making.

**Data breach** – The loss, destruction, alteration, acquisition, or disclosure of information caused by accidental or intentional, unlawful or otherwise unauthorised purposes, which compromise the confidentiality, integrity or availability of information.

**Data consumer** – A person or organisation that uses data to make decisions, take actions, or increase awareness.

**Data processing** – Any operation or set of operations which is performed on data or on sets of data, whether or not by automated means, such as collecting, registering, storing, adapting or altering, cleaning, filing, retrieving, using, disseminating, transferring and retaining or destroying.

**Data protection assessments** – Assessments that aim to identify, evaluate and address the risks to personal data arising from a project, policy, programme or other initiative.

**Data responsibility** – A set of principles, processes and tools that support the safe, ethical and effective management of data in humanitarian response.

**Data security** – A set of physical, technological and procedural measures that safeguards the confidentiality, integrity and availability of data and prevent its accidental or intentional, unlawful or otherwise unauthorised loss, destruction, alteration, acquisition or disclosure.

**Data sensitivity** (distinction) – The Centre for Humanitarian Data distinguishes the following:

> **Sensitive:** Any dataset containing personal data of affected populations or aid workers. Datasets containing demographically identifiable information (DII) or community-identifiable

information (CII) that can put affected populations or aid workers at risk, are also considered sensitive data. Depending on the context, satellite imagery can also fall into this category of sensitivity.

**Uncertain sensitivity**: For this data, sensitivity depends on a number of factors, including other datasets collected in the same context, the technology which is or could be used to extract insights, and the local context from which the data is collected or which will be impacted by the use of the data.

**Non-sensitive**: This includes datasets containing country statistics, roadmaps, weather data and other data with no foreseeable risk associated with sharing.

**Data-sharing paradigms** – The main agreed-upon data-sharing paradigms are (1) limited release, (2) remote access, (3) APIs and open algorithms, (4) precomputed indicators and synthetic data, (5) data collaboratives or spaces. The categorisation is based on the Data-Pop Alliance (2019), Sharing is Caring.

**Data source** – The original collector of data.

**Data transfer** – The act of transferring data or making it accessible to a partner using any means, such as hard copy, electronic means or the internet.

**Federated learning** – A machine-learning technique that trains an algorithm across multiple decentralised devices holding local data without exchanging data between each other. This method enables multiple parties to build a robust machine-learning model in which only differential increment of trained models are transferred among the parties instead of the actual data.

**Harm** – Negative implications of a data-processing initiative on the rights of a data subject, or a group of data subjects, including but not limited to physical and psychological harm, discrimination and denial of access to services.

**Humanitarian data** – Humanitarian data is understood as: (1) data about the context in which a humanitarian crisis is occurring (e.g., baseline or development data, damage assessments, geo-spatial data); (2) data about the people affected by a crisis and their needs; and (3) data about the response by organisations and people seeking to help.

**Humanitarian data ecosystem** – A combined, dynamic overview of data processing activities, data flows and actors interacting with humanitarian data relating to a crisis or geographic area.

**Humanitarian data sharing** – While in the context of this paper, there is a clear distinction between "granting safe access" and "publishing", "sharing" or "releasing" datasets or aggregated insights, for ease of reading, the term is used as an overarching concept that includes all modalities, even if data as such is not shared. In cases where the distinction is critical, the term is specified more closely.

**Identifiers / identifying attributes of data or a record** (distinction) **–** The attributes of a record can be separated into the following categories, which need to be considered in terms of identity disclosure and data and privacy protection.

**Direct identifiers:** Attributes such as name, email or ID that directly lead to an individual's identification without requiring any other information.

**Quasi-identifiers or indirect identifiers**: Attributes such as address, date of birth or job title that may lead to an individual's identification with the help of other information available from other sources.

**Internally displaced persons (IDPs)** – Persons or groups who have been forced or obliged to flee or to leave their homes or places of habitual residence, in particular as a result of or in order to avoid the effects of armed conflict, situations of generalised violence, violations of human rights or natural or human-made disasters, and who have not crossed an internationally recognised state border.

**Microdata** – A set of records containing information on an individual or respondents, households, an establishment or economic entities. Such records are collected through surveys, a census, administrative forms or registries.

**Personal data** – Any information relating to an identified individual, in particular by reference to an identifier such as a name, an identification number, location data, an online identifier or to one or more factors specific to the physical, physiological, genetic, mental, economic, cultural or social identity of that natural person.

**Personally identifiable information (PII)** – Legally described as 'personal data' in Europe and 'personal information' in some other jurisdictions; generally understood as anything that can directly identify an individual, although it is important to note that PII exists along a spectrum of both identifiability and sensitivity.

**Perturbative vs. non-perturbative techniques (in SDC)** – The choice between non-perturbative and perturbative anonymisation methods (in statistical disclosure control), can limit the types of analyses that can be performed on the anonymised version of a dataset without access to additional information on how the original data has been altered.

**Non-perturbative methods** are anonymisation methods that reduce the detail in the data or suppress certain values (masking) without distorting the data structure.

**Perturbative methods** do not suppress values in the dataset but perturb (i.e., alter) values to reduce the risk of disclosure by creating uncertainty around the true values. As no values are suppressed, information loss is reduced. The structure of the data can be altered, however, which has implications that those aiming to use the data for analysis need to be aware of and have the necessary tools to address. Post-randomisation (PRAM), noise addition, shuffling and rank swapping are examples of perturbative methods.

**Profiling (in the displacement context)** – Profiling is a collaborative process for analysing displacement situations. It entails gathering information on populations affected by displacement and the local population in order to take action to advocate on their behalf, to protect and assist them and to help bring about a solution to their displacement. Sensitive individual-level data is collected, which, when disclosed or re-identified, could put individuals at risk. This includes, for instance, (1) the number of

displaced persons, disaggregated by age and sex, (2) location(s), (3) cause(s) of displacement, (4) patterns of displacement, (5) protection concerns, and (6) humanitarian needs.

**Pseudonymisation** – Distinct from anonymisation, pseudonymisation refers to the processing of personal data in such a manner that the personal data can no longer be attributed to a specific individual without the use of additional information, provided that such additional information is stored separately.

**Re-identification** – A process by which de-identified (anonymised) data becomes re-identifiable and can therefore be traced back or linked to one or more individuals or groups through reasonably available means at the time of data re-identification.

**Risk–benefit assessment** – A process for identifying and balancing the benefits and risks related to the processing of data, as well as the likelihood, magnitude and severity of harms that can result from the identified risks in a particular context.

**Risk mitigation** – A process for applying specific measures to prevent or minimise the likelihood of potential risks related to the processing of data, as well as to prevent the occurrence of harms or otherwise minimise their magnitude and severity.

**sdcMicro –** Connected to statistical disclosure control, sdcMicro is an R-based package authored by Templ, M., Kowarik, A. and Meindl, B. with tools for the anonymisation of microdata, i.e., for the creation of public and scientific use files.

**Sensitive data** – Personal data which, if disclosed, may result in discrimination against or repression of the individual concerned, or create a negative impact on an organisation's capacity to carry out activities and on the public perceptions of that organisation. Typically, data relating to health, race or ethnicity, religious/political/armed group affiliation, or genetic and biometric data are considered to be sensitive data.

**Sensitive information** – Attributes such as salary, positive or negative test results, disease or preferences that can be used to manipulate an individual or may lead to different behaviour upon the individual.

**Statistical disclosure control (SDC)** – A set of methods for reducing the risk of disclosing information on individuals, respondents, businesses or other organisations. Such methods are only related to dissemination and are usually based on restricting the amount of or modifying the data released.

# Bibliography

Apfelbeck, F. (2020, November). Evaluation of Privacy-Preserving Technologies for Machine Learning. Medium. Retrieved January 29, 2020, from https://medium.com/outlier-ventures-io/evaluation-of-privacy-preserving-technologies-for-machine-learning-8d2e3c87828c.

APL, Johns Hopkins University, Applied Physics Laboratory, (2018). 2018 Annual Report. Retrieved January 29, 2020, from https://www.jhuapl.edu/Content/documents/2018_Annual_Report.pdf

APL, Johns Hopkins University, Applied Physics Laboratory. Available at https://www.jhuapl.edu/

Benschop, T., and Welch, M. (n.d.). Statistical Disclosure Control for Microdata: A Practice Guide for sdcMicro — SDC Practice Guide documentation. World Bank. Retrieved January 29, 2020, from https://sdcpractice.readthedocs.io/en/latest/index.html.

Benschop, T., and Welch, M. (n.d.).. Statistical Disclosure Control for Microdata: Theory — SDC Theory Guide documentation. World Bank. Retrieved January 29, 2020, from https://sdctheory.readthedocs.io/en/latest/

Centre for Humanitarian Data. (2019). Guidance Note: Statistical Disclosure Control. Retrieved January 29, 2020, from https://centre.humdata.org/wp-content/uploads/2019/07/guidance_note_sdc.pdf.

Data-Pop Alliance. About the OPAL Project. Available at https://datapopalliance.org/opal/.

Data-Pop Alliance. Data-Pop Alliance. Available at https://datapopalliance.org/.

Dwork, C., and Roth, A. (2013). The Algorithmic Foundations of Differential Privacy. Foundations and Trends® in Theoretical Computer Science. 9(3–4): pp. 211–407.

Encrypt, S. (2018). 7 Principles of Privacy by Design. Medium. Retrieved January 29, 2020, from https://medium.com/searchencrypt/7-principles-of-privacy-by-design-8a0f16d1f9ce.

European Commission. (2018). UN OCHA Humanitarian Data Exchange. Retrieved January 29, 2020, from https://ec.europa.eu/knowledge4policy/online-resource/un-ocha-humanitarian-data-exchange_en.

Evans G., King G., Schwenzfeier M., et al. (2020). Statistically Valid Inferences from Privacy Protected Data: pp. 35.

Eyupoglu, C., Aydin, M., Zaim, A., et al. (2018). An Efficient Big Data Anonymization Algorithm Based on Chaos and Perturbation Techniques. Entropy. 20(5): pp. 373.

Farias, A.M. (2019). Private AI – Federated Learning with PySyft and PyTorch. Retrieved January 29, 2020, from https://towardsdatascience.com/private-ai-federated-learning-with-pysyft-and-pytorch-954a9e4a4d4e.

Flowminder. FlowKit (GitHub Documentation). Available at https://flowkit.xyz/.

Flowminder. FlowKit: Unlocking the Power of Mobile Data. Available at
https://digitalimpactalliance.org/wp-
content/uploads/2019/02/FlowKit_UnlockingthePowerofMobileData.pdf

Galtier, M., and Marini, C. (2019). Substra: A framework for privacy-preserving, traceable and
collaborative Machine Learning. France: Substra Foundation. Github. Retrieved January 29,
2020, from https://github.com/SubstraFoundation/welcome

Hu, X., Yuan, M., Yao, J., et al. (2015). Differential privacy in telco big data platform. Proceedings of
the VLDB Endowment. 8(12): pp. 1692–1703.

ICRC and VUB. 2017. Handbook on Data Protection in Humanitarian Action. Available at
https://shop.icrc.org/handbook-on-data-protection-in-humanitarian-action.html.

Ippolito, PP. (2019, July). AI Differential Privacy and Federated Learning. Towards Data Science.
Retrieved January 29, 2020, from https://towardsdatascience.com/ai-differential-privacy-and-
federated-learning-523146d46b85.

JIPS. (2020). Anonymisation of Household Survey Microdata. Project Report.

Letouzé, E., and Oliver, N. (2019). Sharing is Caring: Four Key Requirements for Sustainable Private
Data Sharing and Use for Public Good. NYC: Data Pop Alliance.

Letouzé, E., Pentland, A., and Data-Pop Alliance. (2018). Towards a Human Artificial Intelligence for
Human Development. (2): pp.8.

Li, N., Li, T., Venkatasubramanian, S., et al. (2001). t-Closeness: Privacy Beyond k-Anonymity and -
Diversity. Center for Education and Research Information Assurance and Security, Purdue
University, 10.

Lomas, N. (2019, July). Researchers spotlight the lie of 'anonymous' data. TechCrunch. Retrieved
January 29, 2020, from http://social.techcrunch.com/2019/07/24/researchers-spotlight-the-
lie-of-anonymous-data/

Machanavajjhala, A., Gehrke, J., Kifer, D., et al. (2007). ℓ-Diversity: Privacy Beyond k-Anonymity.
*ACM Transactions on Knowledge Discovery from Data*, 1(1), 12.
https://doi.org/10.1145/1217299.1217302

Madianou, M. (2019). Technocolonialism: Digital Innovation and Data Practices in the Humanitarian
Response to Refugee Crises. *Journal Sage Social Media + Society*, 5(3),
https://doi.org/10.1177/2056305119863146

McDonald, S. (2019). From Space to Supply Chains: A Plan for Humanitarian Data Governance.
*Rochester, NY: Social Science Research Network*. Available at
https://papers.ssrn.com/abstract=3436179.

McMahan H.B., Moore E., Ramage D., et al. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. arXiv:1602.05629 [cs]. : . Available at http://arxiv.org/abs/1602.05629.

De Montjoye, Y.-A., Gambs, S., Blondel ,V., et al. (2018). On the privacy-conscientious use of mobile phone data. *Scientific Data*, 5(1), 180-286, https://www.nature.com/articles/sdata2018286

Montjoye, Y.A., de Gambs, S., Blondel, V., et al. (2018). On the privacy-conscientious use of mobile phone data. *Scientific Data*. 5(1), 1–6. https://www.nature.com/articles/sdata2018286

Narayanan, A., and Shmatikov, V. (2008). Robust De-anonymization of Large Sparse Datasets. 2008 IEEE Symposium on Security and Privacy, 111–125, https://ieeexplore.ieee.org/document/4531148

OPAL. About OPAL. Available at https://www.opalproject.org/about-opal.

OPAL. OPAL Project A Closer Look. Available at https://www.opalproject.org/a-closer-look.

Pinot, R. (2018). Minimum spanning tree release under differential privacy constraints. arXiv:1801.06423 [cs, math, stat]: Available at http://arxiv.org/abs/1801.06423.

Sedlák, V. (2019). De-identification for data journalists. DataJournalism.com. Retrieved January 29, 2020, from https://datajournalism.com/read/longreads/de-identification-for-data-journalists

Sharbain, R. (2019). Data protection policy void threatens privacy rights of citizens and refugees in Jordan. Global Voices Advocacy. Retrieved January 29, 2020, from https://advox.globalvoices.org/2019/12/30/data-protection-policy-void-threatens-privacy-rights-of-citizens-and-refugees-in-jordan/.

Rocher, L., Hendrickx J.M., and de Montjoye Y.A. (2019). Estimating the success of re-identifications in incomplete datasets using generative models. *Nature Communications*. 10(1), 1–9, https://doi.org/10.1038/s41467-019-10933-3

UN OCHA. Frequently Asked Questions – Humanitarian Data Exchange. Available at https://data.humdata.org/faq.

UN OCHA. Humanitarian Data Exchange. Available at https://data.humdata.org/.

UN OCHA. The Centre for Humanitarian Data (HDC). Available at https://centre.humdata.org/.

UN OCHA. UN-OCHA Humanitarian Data Exchange Project GitHub Documentation. Available at https://github.com/OCHA-DAP.


UNHCR. 2018. UNHCR 2018 Global Strategic Priorities Progress Report. Geneva. Available at http://reporting.unhcr.org/sites/default/files/2018%20Global%20Strategic%20Priorities%20Progress%20Report.pdf.

Yu, S. (2016). Big Privacy: Challenges and Opportunities of Privacy Study in the Age of Big Data. IEEE Access. 4. pp.2751–2763.

Zumel, A.N. (2015, October). A Simpler Explanation of Differential Privacy. Win-Vector Blog Retrieved January 29, 2020, from http://www.win-vector.com/blog/2015/10/a-simpler-explanation-of-differential-privacy/

*You see things; and you say "Why?"*
*But I dream things that never were;*
*and I say "Why not?"*

— George Bernard Shaw

*Lovely!*
*It also fits to causal inference*
*and counter-factual questions!*

— Navid Rekabsaz