

Differential privacy for public health data

An innovative tool to optimize information sharing while protecting data confidentiality

Dyda, Amalie; Purcell, Michael; Curtis, Stephanie; Field, Emma; Pillai, Priyanka; Ricardo, Kieran; Weng, Haotian; Moore, Jessica C.; Hewett, Michael; Williams, Graham; Lau, Colleen L.

Published in:
Patterns

DOI:
[10.1016/j.patter.2021.100366](https://doi.org/10.1016/j.patter.2021.100366)

Published: 10/12/2021

Document Version
Publisher's PDF, also known as Version of record

[Link to publication](#)

Citation for published version (APA):

Dyda, A., Purcell, M., Curtis, S., Field, E., Pillai, P., Ricardo, K., Weng, H., Moore, J. C., Hewett, M., Williams, G., & Lau, C. L. (2021). Differential privacy for public health data: An innovative tool to optimize information sharing while protecting data confidentiality. *Patterns*, 2(12), 1-7. [100366].
<https://doi.org/10.1016/j.patter.2021.100366>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Perspective

Differential privacy for public health data: An innovative tool to optimize information sharing while protecting data confidentiality

Amalie Dyda,^{1,9,*} Michael Purcell,^{2,9} Stephanie Curtis,³ Emma Field,^{4,5} Priyanka Pillai,⁶ Kieran Ricardo,² Haotian Weng,² Jessica C. Moore,⁷ Michael Hewett,⁸ Graham Williams,² and Colleen L. Lau^{1,4}

¹School of Public Health, University of Queensland, 288 Herston Road, Herston, QLD 4006, Australia

²Software Innovation Institute, Australian National University, CSIT Building (#108), North Road, Acton, ACT 2601, Australia

³National Centre for Epidemiology and Population Health, Australian National University, Acton, ACT, Australia

⁴Research School of Population Health, Australian National University, 62 Mills Road, Acton, ACT 2601, Australia

⁵Global and Tropical Diseases Division, Menzies School of Health Research, Charles Darwin University, Darwin, NT, Australia

⁶Doherty Institute, 792 Elizabeth Street, Melbourne, VIC 3000, Australia

⁷School of Computing, Australian National University, CSIT Building (#108), North Road, Acton, ACT 2601, Australia

⁸The National Centre for Geographic Resources & Analysis in Primary Health Care (GRAPHC), Research School of Population Health, Australian National University, 62 Mills Road, Acton, ACT 2601, Australia

⁹These authors contributed equally

*Correspondence: a.dyda@uq.edu.au

<https://doi.org/10.1016/j.patter.2021.100366>

THE BIGGER PICTURE Differential privacy is an innovative technique that can be applied to data to protect confidentiality. This has been used primarily to protect private sector data, but has significant implications for public health. We describe the methods of differential privacy in terms understandable to a non-computer-science audience. To our knowledge, this is the first article describing differential privacy in language and context appropriate for a health audience. The case study described shows the feasibility of the use of differential privacy for public health surveillance data to optimize information sharing while protecting data confidentiality. This method allows for data to be released in more granular detail in terms of time, place, and person without compromising privacy and confidentiality. Future research needs to consider other use cases, including a range of surveillance systems and applications in other types of health data.



Concept: Basic principles of a new data science output observed and reported

SUMMARY

Coronavirus disease 2019 (COVID-19) has highlighted the need for the timely collection and sharing of public health data. It is important that data sharing is balanced with protecting confidentiality. Here we discuss an innovative mechanism to protect health data, called differential privacy. Differential privacy is a mathematically rigorous definition of privacy that aims to protect against all possible adversaries. In layperson's terms, statistical noise is applied to the data so that overall patterns can be described, but data on individuals are unlikely to be extracted. One of the first use cases for health data in Australia is the development of the COVID-19 Real-Time Information System for Preparedness and Epidemic Response (CRISPER), which provides proof of concept for the use of this technology in the health sector. If successful, this will benefit future sharing of public health data.

BACKGROUND

The coronavirus disease 2019 (COVID-19) pandemic has highlighted the need for the timely collection and sharing of public health data to provide information for policy makers and frontline workers to make rapid and informed decisions. It is important that sharing of data is balanced with protecting the confidentiality of an individual's health information. In this paper we

discuss a cybersecurity method known as differential privacy, an innovative mechanism that could be used to optimize information sharing while protecting the confidentiality of public health surveillance data. Privacy refers to an individual's right to decide whether information about him or her is released, while confidentiality is an assurance given by a data holder that they will not violate any individual's privacy by releasing data the individual desires to be private. Privacy prevents information about a



person being shared, and confidentiality ensures data relating to that individual are shared only with authorized parties, while his or her identity is protected.¹

PREVIOUS BREACHES OF HEALTH DATA CONFIDENTIALITY IN AUSTRALIA

Concerns regarding data confidentiality are warranted, with numerous previous data breaches occurring in relation to health data in Australia. In one such example, a graph from South Australia Health reporting data on children treated at hospital for respiratory infection, gastroenteritis, and whooping cough was displayed online from 2005 to 2018. The graph was linked to the source data, including names, date of birth, and test results, and had recorded 300 views before the breach was identified and the data were removed.² Another example of a serious breach involved data provided by individuals to the Red Cross Blood Service between 2010 and 2016. Personal and health information, including details about “at-risk behavior,” were completed in an online application form. The file containing all these data was moved to an unsecure device and accessed by an unauthorized party.³ Both of these examples involve datasets that contained personal identifiers.

There are also instances in which datasets have been shared with measures put in place to protect confidentiality, such as the removal of primary identifiers (e.g., name), but breaches have occurred. For example, on August 1, 2016, the Australian Department of Health publicly released medical billing records of 2,985,511 unique individuals on the data.gov.au website. The records included Medicare Benefits Scheme (MBS) data from 1984 to 2014 and Pharmaceutical Benefits Scheme (PBS) data from 2003 to 2014, containing historical health data of around 10% of the population, including details on services provided by doctors, pathologists, diagnostic imaging, and allied health.⁴ The dataset was released for research in the public interest and used anonymized unique identifiers while applying numerous confidentiality measures to prevent information from being identifiable, including encryption, perturbation, aggregation, and exclusion of rare events. The dataset was downloaded approximately 1,500 times while it was publicly available.

In December 2016, the Australian Department of Health was alerted by computer science researchers at the University of Melbourne that it was possible to decrypt Medicare service provider unique identifiers.⁴ The University of Melbourne researchers used cryptographic attack methods (finding a security weakness, usually in the code) to recover information from the dataset. The researchers did not reveal personal information when reporting their findings, so no individual’s privacy was compromised. Rather, it was used as an example of the importance of data security and the challenges of performing de-identification. The authors outlined their methods in a paper, “Health Data in an Open World,”⁴ which explains how individuals could be re-identified through decryption and linking.

HOW PREVIOUS DATA BREACHES HAVE AFFECTED ACCESS TO SURVEILLANCE DATA IN AUSTRALIA

Following an investigation into the MBS and PBS data breaches described above, it was concluded that decision-making at the

Australian Department of Health for releasing these data was not a clear and documented approval process, and there were no rigorous risk management processes or any significant degree of cross-government coordination.⁵ Consequently, the Australian government published guidance on the “Process for Publishing Sensitive Unit Record Level Public Data as Open Data,”⁶ which provides guidance on releasing datasets related to sensitive information. In addition, the Commonwealth Government Privacy Amendment (Re-identification Offence) Bill 2016 was introduced as a method of deterrence to make it an offense for an entity to intentionally re-identify information that has been disclosed by a Commonwealth agency.⁷ Overall, there was an emphasis on the failure of the government, rather than appreciation for the relationship between researcher and agency, to ensure privacy breaches did not occur. The Notifiable Data Breaches program was also introduced in 2018, which requires that organizations notify individuals and the Office of the Australian Information Commissioner when a data breach, which may cause serious harm, occurs.⁸

CURRENT STRATEGIES TO ENSURE DATA PRIVACY

When data custodians allow access to data, several techniques can be used to protect privacy, which focus on altering either datasets or processes around data access, including the following:

- **De-identification** is the process of (1) removing personal identifiers (such as name and address) and typically replacing them with token identifiers and (2) removing or altering other information that may permit the re-identification of an individual through linkage with other data whereby persons are identified.⁹ This process should also consider whether data are potentially identifiable, even after the name and address are removed. Data, particularly small datasets, can be potentially re-identifiable through indirect identification, using a combination of variables such as date of birth and postcode, or these variables can be used to significantly narrow down the number of people who would fit the set of criteria.^{10,11}
- **Aggregation** involves presenting data in a summarized format, particularly when there is a small number of people reported in a geographic area or as some other subgroup (e.g., by ethnicity).¹²
- **Authentication** is the process of identifying the person(s) accessing the data. Digital identity verification is a common feature of authentication, which results in having multiple “digital identities” through the use of different log-in and password combinations.¹³
- **Authorization** is providing defined roles or operations for how the data are accessed and used.¹⁴
- **Encryption** is the process of using an algorithm to scramble data, and a key is used by those with authorization to unscramble and decrypt the data.¹³

A POTENTIAL INNOVATIVE SOLUTION: DIFFERENTIAL PRIVACY

A possible solution to improve data security when making sensitive health data securely available is the use of the recently

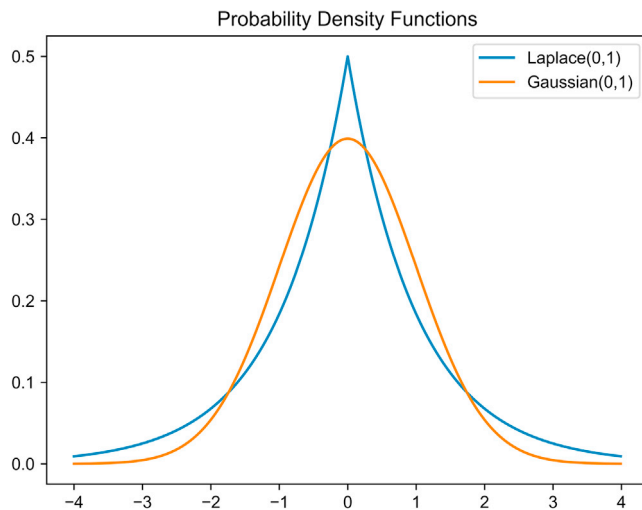


Figure 1. Comparative plot of the density functions of the Laplace (0; 1) and the Gaussian (0; 1) distributions

Note that the Laplace distribution has a sharp “peak” at zero, while the Gaussian is more rounded. Also note that the tails of the Laplace distribution are much heavier than those for the Gaussian distribution. That is, samples drawn from the Laplace distribution are more likely to be farther away from the mean than are samples drawn from the Gaussian distribution.

developed cybersecurity method known as differential privacy. Differential privacy is a mathematically rigorous definition of privacy that aims to protect against all possible adversaries who might want to compromise privacy and confidentiality. It does this by limiting the information gained by the worst possible adversary: someone who knows all but one row of the database, has infinite computing power, and attempts to discover the remaining row of the database through targeted queries. In layperson’s terms, statistical noise is applied to the data so that overall epidemiological patterns can be described but individual data cannot be extracted. To date, this has primarily been used in the private sector and in the area of data science, with application of the method by Apple and Google,^{15,16} as well as in the release of the 2020 US census data.¹⁷ The application of differential privacy for health data is starting to be recognized, with some researchers in the field recommending its use for creating health datasets for secondary purposes,¹⁸ and one pilot study applying the technique to query a patient database to select individuals for recruitment into a clinical trial.¹⁹

CASE STUDY

What would a worst-case adversary look like? Imagine a scenario where the adversary, Eve, knows everything about everyone in her community except Alice. Suppose also that Eve has unlimited computational resources. Eve’s goal is to analyze any data that have been released about the people in her community to try to determine if Alice has been diagnosed with COVID-19.

- Is it safe to release aggregate data? No. If Eve has access to any exact statistic, like the mean, count, or variance, that included Alice, she can easily compute what these statistics would be in the case where Alice has/has not been

diagnosed with COVID-19. Eve can then compare these to the true values and deduce whether Alice has been diagnosed with COVID-19.

- Will low-count cells suppression work? No, the privacy attack detailed above works regardless of how many people are in the community.
- Will methods like binning work? Binning is a method that allows the grouping of continuous values into a smaller number of groups or “bins.”²⁰ This may be useful for some types of health data but not for all. Binning will sometimes protect Alice’s identity from Eve, but is much less likely to be effective at protecting everyone’s identity from all possible attackers. Suppose:
 - the bins are counts of COVID-19 notifications of 0–5 and 6–10 cases,
 - the true number of COVID-19 cases is 6, and
 - Alice has been diagnosed with COVID-19.

Eve uses her ancillary data to calculate that the number of COVID-19 cases, excluding Alice, is 5. Eve compares this to the released data that state that the number of cases is in the 6–10 bin and deduces that Alice has been diagnosed with COVID-19.

One of the most common ways to achieve differential privacy is to add Laplace noise to the statistics.²¹ To do so, we first compute the exact statistics, then generate random values distributed according to the Laplace distribution, and finally add the random (noise) values to the exact statistics. This results in a differentially private version of the original data, which are the statistics that are released. We have used the term statistics here, given that this is common terminology in public health. However, in the language of differential privacy, the data that we release would more commonly be called the “query responses” and will be referred to as this henceforth.

The Laplace distribution is similar to the more familiar Gaussian (normal) distribution (as shown in Figure 1) in that both are symmetric about some mean value and that random samples drawn from either distribution will tend to be tightly clustered around the mean. The Laplace distribution has somewhat heavier tails than the Gaussian distribution, however, and is more likely to produce samples that are far from its mean. It turns out that this property is critical for ensuring that the addition of Laplace noise is sufficient to guarantee differential privacy.

Following the case study above, assume that:

- Eve knows the COVID-19 status of all community residents except Alice;
- Eve knows that there are 999 cases in her community. That is, if Alice has not been diagnosed with COVID-19, there will be 999 cases. Or if Alice has been diagnosed, there will be 1,000 cases;
- random Laplace noise is added to the count before being released to Eve;
- the query response released to Eve is 998;
- Eve knows the privacy protocol;
- Eve believes both outcomes are equally likely. That is, before observing the query response, Eve believes that there is a 50% chance that Alice has been diagnosed with COVID-19.

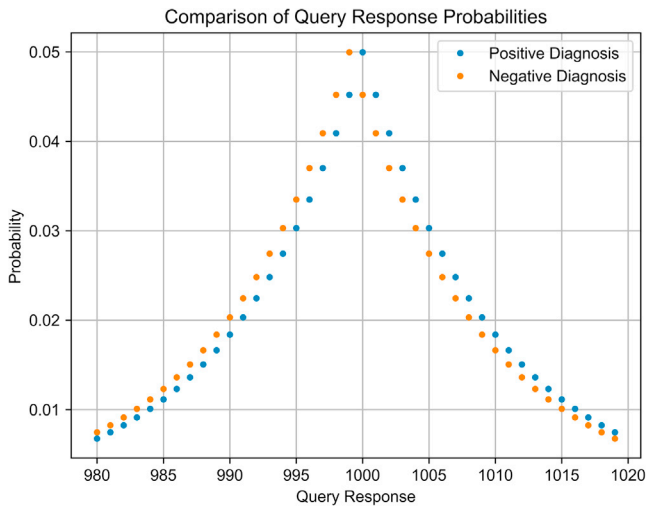


Figure 2. Distribution of the probabilities of query responses produced by the Laplace mechanism when Alice has been diagnosed with COVID-19 (blue) and when Alice has not been diagnosed with COVID-19 (orange).

Figure 2 shows the distribution of the values produced by the Laplace mechanism when Alice has or has not been diagnosed with COVID-19. The distributions of the probabilities of query responses under the two scenarios are virtually indistinguishable.

From Eve’s ancillary data, she knows that either there were 999 cases and Alice has not been diagnosed with COVID-19, or there were 1,000 cases and Alice has been diagnosed with COVID-19. Because Eve knows the privacy protocol, she can calculate the probabilities that any count would be released under each scenario. Because Eve believes both outcomes are equally likely, Eve’s optimal strategy is to guess that Alice’s diagnosis is whatever maximizes the likelihood of the observed query response. For example, suppose that the query response is 998. Eve knows that the probability of this happening is 4.1% if Alice has been diagnosed with COVID-19 and 4.5% otherwise. So, Eve’s best guess is that Alice has not been diagnosed with COVID-19. Suppose, however, that the query response is 1,003. Eve knows that the probability of this happening is 3.8% if Alice has been diagnosed with COVID-19 and 3.3% otherwise. So, in this case Eve’s best guess is that Alice been diagnosed with COVID-19. Importantly, regardless of the actual query response, we can apply Bayes formula (a mathematical formula used for calculating conditional probabilities) to show that there is a 47.5% chance that Eve’s guess is incorrect. Therefore, the best Eve can do is little better than random guessing.

HISTOGRAM QUERIES

One particularly important class of queries that can be made differentially private by addition of Laplace noise is histogram queries. A histogram query returns the set of values for the rows in the database and the number of rows that have each possible value. For example, consider a database that consists of the age group of every person who had a COVID-19 test on a given date. A histogram query on that database simply counts the number of persons in each age group.

Histogram queries can be made differentially private by adding independent Laplace noise to the count for each possible value. Somewhat surprisingly, the scale of the noise that must be added to achieve a given level of privacy (see “Privacy budget” below) does not depend on the number of possible values. For example, the amount of noise that must be added to each count is the same whether there are 5 possible values or 500. Furthermore, the scale of the noise does not depend on the population size. Therefore, if the population size is large and the number of possible values is small, then the magnitude of the noise will be small relative to the true counts.

Returning to the example COVID-19 database described above, suppose that the persons’ ages were reported in eight categories: 0–9, 10–19, 20–29, 30–39, 40–49, 50–59, 60–69, and 70+ years. A histogram query on this database would then return eight counts, one for each possible age group. This query can be made differentially private by adding independent Laplace noise to each of these counts. Figure 3 depicts the results of a non-private histogram query alongside a differentially private (perturbed counts) version of the same query, showing very similar overall patterns of age distribution.

PRIVACY BUDGET

The strength of the privacy guarantee provided by a differentially private release mechanism depends on the scale of the Laplace noise that is added. In general, adding more noise yields stronger privacy. The strength of this privacy guarantee, or the privacy budget, is based on the value of a privacy parameter called epsilon (ϵ). Smaller values of epsilon yield stronger privacy guarantees.^{22,23} Unfortunately, smaller values of epsilon also require adding more noise to the exact query responses, making them less accurate. When releasing data, an appropriate balance between accuracy and privacy needs to be determined. This is known as the privacy-utility trade-off.

The privacy guarantees described above apply only when a differentially private release mechanism is used to answer a single query. In practice, however, a release mechanism may be required to compute responses to many queries. For example, in a dataset containing COVID-19 data, an individual may search for all COVID-19 cases in a specific location, for which a differentially private result would be produced. However, each time this search is run, a new differentially private result would be produced. Each resulting search is considered an additional query. Answering multiple queries weakens the privacy guarantee that any release mechanism can provide. Hence, if a system using differential privacy allows more than one query, there should be a limit on the number of queries allowed. If an unlimited number of queries were allowed, this may make the data re-identifiable. To see why, suppose that the Laplace mechanism is used to compute many independent responses to a query. In this case, an attacker could compute the average of the noisy query responses to obtain a very accurate estimate of the exact query response, which would defeat the purpose of applying noise in the first place! One of the most important properties of differentially private release mechanisms is that they allow us to quantify the strength of the privacy guarantee that a release mechanism can provide when used to answer multiple queries. In general, halving the privacy parameter allows us to answer twice as

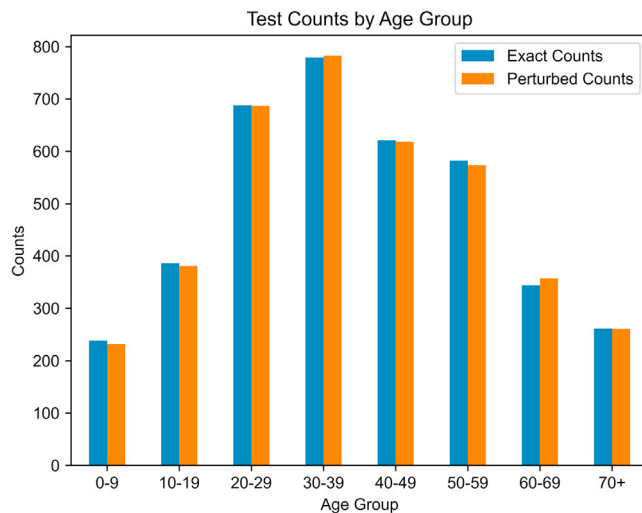


Figure 3. Histogram of real data (blue) compared with differentially private query responses of the same dataset ($\epsilon = 1/8$; orange).

many queries while providing the same privacy guarantee. This property is known as composability.²⁴

Composability offers a way to provide strong privacy guarantees when a differentially private release mechanism is used to respond to multiple queries. To do so, we need to calculate the value of epsilon that we are going to use (i.e., how much noise are we going to add) for each query. To calculate the privacy budget, we compute the value of epsilon for each query individually and then add them up. The total is then compared with a chosen threshold to ensure that the system provides a strong enough privacy guarantee. As with individual queries, systems with smaller privacy budgets provide stronger privacy guarantees. One commonly applied rule of thumb is that epsilon “should be thought of as a small number, between approximately 1/1,000 and 1.”²⁵

PRACTICAL CONSIDERATIONS

One major barrier to the widespread adoption of differential privacy is the paucity of available software tools that provide secure and efficient implementations of differentially private release mechanisms and that can be used by those without computer science expertise. Some of the more popular existing tools^{22,26–28} provide implementations of commonly used differentially private release mechanisms for use in the Python, R, SQL, and C++ programming languages. Unfortunately, even the simplest of differentially private release mechanisms, such as the Laplace mechanism described above, require careful implementation to ensure adequate protection for real-world data.²⁹ Furthermore, it can be difficult for non-experts to effectively wield the tools that are available, since doing so requires an understanding of how to choose an appropriate value for epsilon and of how to compute the sensitivity of a given query. To address these shortcomings, one of the authors (M.P.) has developed the Python library ReIM,³⁰ which provides easy-to-use, secure implementations of many fundamental differentially private release mechanisms. The package includes a synthetic

dataset designed to reflect possible real-world COVID-19 testing data and examples of how ReIM can be used to release differentially private histograms as per the example described above.

ADVANTAGES AND LIMITATIONS OF USING DIFFERENTIAL PRIVACY FOR HEALTH DATA

There are a number of distinct advantages to the use of differential privacy for health data. The examples described previously are of small-scale datasets with minimal variables to explain the concept. In real-world public health datasets, however, there are often thousands of counts of disease incidence with a large number of variables, such as age, sex, ethnicity, and geographical location. We suggest that differential privacy in this context can be even more useful. The rigorous privacy guarantees that differential privacy provides make it possible to determine how many counts can be safely released given some fixed privacy budget. For sufficiently large populations it is possible to release thousands of differentially private query responses with both small relative errors and strong privacy guarantees. So, it can be possible to use differential privacy to release query responses in forms such as dashboards that allow for interrogation of the data across numerous variables.

In addition to providing added security, the primary datasets can remain with the data custodian so that it is possible for users to make queries even though the data have not been specifically released to those conducting analyses, e.g., public health practitioners and researchers. Data privacy can also be increased by limiting the number of queries from each user. Differential privacy is most suitable for large datasets for the presentation of counts, rates, and summary statistics (with a caveat that users accept that the released data have been perturbed).

There are also limitations to the use of this technology. Differential privacy is less suitable for data where there are low counts, as the amount of noise applied will more significantly affect the results. For example, reporting a differentially private estimate of 300 cases of COVID-19 in a postcode compared with the true value of 297 has minimal impact on the usefulness of the data. However, reporting a differentially private estimate of eight cases compared with the true value of four has a larger impact on accuracy and usefulness. In instances with very low counts, differential privacy may also result in negative counts; hence, further work is needed to extend differential privacy for very low value counts to be used in public health surveillance. Similarly, differential privacy is not useful in instances where exact counts are required, or when exact counts are released by other sources, e.g., media reports of the exact number of new COVID-19 cases.

POSSIBLE USES FOR PUBLIC HEALTH

With the added protection afforded by differential privacy, it may be possible to share more detailed information about health and disease patterns in a public forum. This would allow clinicians, public health practitioners, researchers, and the general public to access more timely and detailed information about diseases. This is particularly important during rapidly evolving situations such as public health emergencies like outbreaks and pandemics. For example, during the COVID-19 pandemic, a number

of inquiries have identified the need to do better at including general practitioners and other primary health care workers in the response, including improving their ability to access detailed data about disease epidemiology and spread.^{31,32} Currently, releasing details about COVID-19 cases such as exposure site and demographics (e.g., indigenous status and age) and post-code could potentially enable identification of individuals who have tested positive for COVID-19. Differential privacy could significantly improve protection of this type of data release.

One of the first uses of differential privacy for health purposes in Australia is the development of the COVID-19 Real-Time Information System for Preparedness and Epidemic Response (CRISPER) by our team. This system includes multiple dashboards to visualize and interact with data, including data on confirmed cases and deaths, source of infection, contact tracing alert locations, and laboratory testing. CRISPER's interactive mapping tool allows users to interrogate data based on time *and* place *and* source of infection, and answer questions specific to their information needs.³³ At the time of writing, these functionalities were not available through health department websites. For example, general practitioners could obtain information about the number of locally acquired new cases over the past 14 days in the postcodes surrounding the clinic, and use this information for risk assessment for themselves, their staff, and their patients. CRISPER is designed to use a differential privacy algorithm through a data engine to protect data that are not publicly available (e.g., data stratified by age, sex, and comorbidities). Access to line-listed data has been the major challenge for the CRISPER project, and the system currently uses data scrapers and application programming interfaces (APIs) to parse data from a number of public sources such as health department websites. However, once access to line-listed data has been approved, the system is ready to ingest more detailed data and release differentially private query responses. The use of differential privacy in the CRISPER project provides a proof of concept of the application of this innovative mechanism, which may encourage sharing of public health data in the future by providing enhanced privacy and confidentiality.

ACKNOWLEDGMENTS

We would like to acknowledge that funding for this article was provided by a grant from the Australian Partnership for Preparedness Research on Infectious Disease Emergencies (APPRISE)-National Health and Medical Research Council (NHMRC) Centre for Research Excellence. C.L.L. was supported by an NHMRC Fellowship (APP1193826).

AUTHOR CONTRIBUTIONS

All authors developed the original idea and planned the writing of the manuscript. A.D., E.F., P.P., S.C., and C.L.L. contributed to the research underpinning the manuscript. M.P., K.R., H.W., G.W., J.M., and M.H. assisted with the computer science interpretation for a public health audience and provided analysis for examples used in the paper. A.D. and M.P. led the writing of the manuscript. All authors provided ongoing feedback and approved the final version of the manuscript.

DECLARATION OF INTERESTS

The authors declare no competing interests.

REFERENCES

- Geiderman, J.M., Moskop, J.C., and Derse, A.R. (2006). Privacy and confidentiality in emergency medicine: obligations and challenges. *Emerg. Med. Clin. North Am.* 24 (3), 633–656.
- Thomas-Wilson S. 7000 Patient records from Women's and Children's hospital exposed online in embedded data. *Advertiser*. August 4, 2018.
- Australian Red Cross. Blood Service Apologises for Donor Data Leak 28th October 2016. <https://www.donateblood.com.au/media/news/blood-service-apologises-donor-data-leak>.
- Culnane, C., Rubinstein, B.I.P., and Teague, V. (2017). Health data in an open world. *CoRR abs/1712.05627*, arXiv:1712.05627.
- Commissioner AGOotAI (2018). Publication of MBS/PBS Data: Commissioner Initiated Investigation Report.
- Australian Government Department of the Prime Minister and Cabinet. Process for Publishing Sensitive Unit Record Level Public Data as Open Data. https://blog.data.gov.au/assets/files/process_for_publishing_open_data_dec16.pdf.
- Commonwealth Government Privacy Amendment (Re-identification Offence) Bill, Parliament of Australia, 2016. https://www.aph.gov.au/Parliamentary_Business/Bills_LEgislation/Bills_Search_Results/Result?bld=s1047.
- Australian Government Office of the Australian Information Commissioner (2021). Notifiable Data Breaches. <https://www.oaic.gov.au/privacy/notifiable-data-breaches/about-the-notifiable-data-breaches-scheme/>.
- Australian Government Office of the Australian Information Commissioner (2018). De-identification and the Privacy Act March. <https://www.oaic.gov.au/privacy/guidance-and-advice/de-identification-and-the-privacy-act/>.
- Queensland Government Department of Health (2021). What is 'potentially identifiable' data?. <https://www.health.qld.gov.au/hsu/potentially-identifiable-data>.
- Rocher, L., Hendrickx, J., and de Montjoye, Y. (2019). Estimating the success of re-identifications in incomplete datasets using generative models. *Nat. Commun.* 10, 3069.
- Wen, T. (2020). Data aggregation. In *Encyclopedia of Big Data*, L.A. Schintler and C.L. McNeely, eds. (Springer International Publishing), pp. 1–4.
- National Research Council (US) (1997). Committee on Maintaining Privacy and Security in Health Care Applications of the National Information Infrastructure. For the Record Protecting Electronic Health Information (National Academies Press (US)).
- Institute of Medicine (US) (1994). Committee on regional health data networks. Health data in the information age: use, disclosure, and privacy, M.S. Donaldson and K.N. Lohr, eds. (National Academies Press (US)).
- Apple Inc. (2021). Apple Differential Privacy Technical Overview (Apple).
- Erlingsson, Ú., Pihur, V., and Korolova, A. (2014). RAPPOR: randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1054–1067.
- Petti, S., and Flaxman, A.D. (2020). Differential privacy in the 2020 US census: what will it do? Quantifying the accuracy/privacy tradeoff. *Gates Open Res.* 3, 1722.
- Dankar, F., and Emam, K. (2012). The Application of Differential Privacy to Health Data (ACM International Conference Proceeding Series).
- Dankar, F., and Emam, K. (2013). Practicing differential privacy in health care: a review. *Trans. Data Privacy* 6, 35–67.
- Lin, Z., Hewett, M., and Altman, R.B. (2002). Using binning to maintain confidentiality of medical data. *Proc. AMIA Symp.* 454–458.
- Kifer, D., Messing, S., Roth, A., Thakurta, A., and Zhang, D. (2020). Guidelines for Implementing and Auditing Differentially Private Systems (arXiv).

22. Alvim, M., Andrés, M., Chatzikokolakis, K., Degano, P., and Palamidessi, C. (2011). Differential Privacy: On the Trade-Off between Utility and Information Leakage. https://doi.org/10.1007/978-3-642-29420-4_3.
23. Lee, J., and Clifton, C. (2011). How Much Is Enough? Choosing ϵ for Differential Privacy, pp. 325–340. https://doi.org/10.1007/978-3-642-24861-0_22.
24. Dwork, C., and Roth, A. (2014). The algorithmic foundations of differential privacy. *Found Trends Theor. Comput. Sci.* 9, 211–407.
25. Wood, A., Altman, M., Bembenek, A., Bun, M., Gaboardi, M., Honaker, J., Nissim, K., O'Brien, D., Steinke, T., and Vadhan, S. (2018). Differential privacy: a primer for a non-technical audience. *Vanderbilt J. Entertainment Technol. L.* 27, 209.
26. Holohan, N., Braghin, S., Aonghusa, P., and Levacher, K. (2019). Diffprivlib: The IBM differential privacy library. *arXiv*, 1907.02444 [cs.CR].
27. Rubinstein, B.I.P., and Aldà, F. (2017). Pain-free random differential privacy with sensitivity sampling. In *Proceedings of the 34th International Conference on Machine Learning; Proceedings of Machine Learning Research*, P. Doina and T. Yee Whye, eds. (PMLR), pp. 2950–2959.
28. Wilson, R., Zhang, C., Lam, W., Desfontaines, D., Simmons-Marengo, D., and Gipson, B. (2019). Differentially private SQL with bounded user contribution. *arXiv*, 1909.01917 [cs.CR].
29. Mironov, I. (2012). On significance of the least significant bits for differential privacy. In *Proceedings of the 2012 ACM Conference on Computer and Communications Security (Association for Computing Machinery)*, pp. 650–661.
30. Purcell, M., Ricardo, K., and Nabaglo, J. (2021). ReIM. GitHub Repository <https://github.com/antoyo/reim>.
31. Parliament of Victoria Public Accounts and Estimates Committee (2020). Inquiry into the Victorian Government's Response to the COVID-19 Pandemic: Interim Report. <https://www.parliament.vic.gov.au/paec/inquiries/inquiry/1000>.
32. Tsirtsakis, A. (2020). Inquiry told GPs should be 'essential' part of pandemic preparedness. NewsGP <https://www1.racgp.org.au/newsgp/professional/inquiry-told-gps-should-be-essential-part-of-pande>.
33. Field, E., Dyda, A., and Lau, C.L. (2021). COVID-19 real-time information system for preparedness and epidemic response (CRISPER). *Med. J. Aust.* 214, 386–386.e1.

About the author

Dr. Amalie Dyda is an infectious disease epidemiologist working as a senior lecturer at the University of Queensland. She completed a Master of Applied Epidemiology at the Australian National University (ANU) in 2010 and a PhD investigating vaccine-preventable diseases in adults at the University of New South Wales in 2017. She has previously worked in numerous health departments throughout Australia and has expertise in infectious diseases, outbreak investigation, analysis of large datasets, and health informatics.

Dr. Michael Purcell is a mathematician and computer scientist with a background in probability, mathematical statistics, cryptography, and machine learning. He was awarded his PhD from the University of Utah in 2010. Subsequently, he worked as an applied research mathematician for the US government. He currently works as a senior research translation engineer for the Software Innovation Institute, which is itself a part of the ANU College of Engineering and Computer Science. His current research interests include practical applications of differential privacy and secure multiparty computation.