



Please review the Supplemental Files folder to review documents not compiled in the PDF.

HIPAA and the Torrential Leak of 'De-Identified' Electronic Health Record Data

Journal:	<i>New England Journal of Medicine</i>
Manuscript ID	21-02616
Article Type:	Perspective
Date Submitted by the Author:	12-Feb-2021
Complete List of Authors:	Mandl, Kenneth; Boston Children's Hospital, Computational Health Informatics Program; Perakslis, Eric; Duke Clinical Research Institute
Abstract:	

SCHOLARONE™
Manuscripts

HIPAA and the Torrential Leak of ‘De-Identified’ Electronic Health Record Data

Kenneth D. Mandl, MD, MPH

Computational Health Informatics Program, Boston Children’s Hospital, Boston, MA

Department of Biomedical Informatics, Harvard Medical School, Boston, MA

Eric D Perakslis, PhD

Duke Clinical Research Institute, Duke University Medical Center, Durham, NC

Correspondence to:

Kenneth D. Mandl, MD, MPH

Director, Computational Health Informatics Program

Boston Children’s Hospital

300 Longwood Avenue-LM5506, Mailstop BCH3187

Boston, Massachusetts 02115

617.355.4145

kenneth_mandl@harvard.edu

Word count: 1209

1
2
3 The permissible sharing of identified data to be used for treatment, payment and
4 operations, by covered entities with their business associates, has led to a torrent of electronic
5 health record (EHR) data flowing out of healthcare provider silos. The Health Insurance
6 Portability and Accountability Act (HIPAA) also permits business associates to de-identify data
7 on behalf of the covered entity, and once those data are de-identified, the business associate may
8 use them freely, unless contractually prohibited from doing so. These circumstances enabled the
9 rise of a multibillion-dollar industry comprising dozens of health data aggregation companies
10 and hundreds of health data tool and technology companies aggregating, linking, and monetizing
11 EHR data.
12
13
14
15
16
17
18
19
20
21
22
23

24 This phenomenon has been amplified by the explosion of data production after the
25 HITECH Act promoted widespread adoption of EHRs, predicated on the aspiration for accurate
26 and complete patient records, rapid real-time learning, better coordinated care, accelerated
27 biomedical discovery, and digital and data exchange directly with patients. It is ironic that while
28 physicians and their patients still find it difficult to obtain complete medical record information
29 in a timely fashion, a provision in the HIPAA privacy rule permits massive troves of their digital
30 healthcare data to traverse the medical-industrial complex unmonitored and unregulated.
31
32
33
34
35
36
37
38
39

40 While some may question the necessity or viability of privacy in the digital age, privacy
41 is essential for limiting government and industry power, bolstering self-determination and
42 individual preferences, and preserving reputations. Though the HIPAA Privacy Rule governs
43 covered entities' uses of identifiable data, it does not apply when data are considered to be de-
44 identified by expert determination or the safe harbor method, which requires that 18 specified
45 identifiers are removed.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 On the one hand, assembly of vast datasets advances the laudable objectives of a learning
4 healthcare system, whereby the data routinely collected during the care delivery process
5 continuously drive ever more intelligent treatment decisions. However, markets for secondary
6 use of patient data, as configured today, do not always serve the best interests of patients or the
7 public. Consider, for example, a data aggregation company that drives up the price of drugs by
8 targeting physicians and patients for pharmaceutical detailing.
9

10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Though deidentification is often treated as infallible, it is not. Individuals can often be reidentified using only a handful of attributes.¹ Further, de-identification technologies relying on encryption may yet be highly vulnerable to advances in computing. When there are no contractual controls to govern data produced by the covered entity that is shared with a business associate, if something goes wrong, the harms accrue only to the patients; there is no comprehensive data privacy law in the U.S and none of the patchwork of laws or regulations protects patients from potential uses of de-identified data. Absent are duty to report when data has been re-identified or linked to external data sources (e.g., financial records). There is no opportunity for an individual seeking redress for re-identified data, and no external way to verify or ensure adequacy of de-identification.⁴

Research on de-identified data rarely is under institutional review board (IRB) oversight. Further, the free market for de-identified data can preclude desirable rigor on data quality, especially as an analyst or researcher cannot go back to the original data source for removal of duplicates, proper data validation, or correction of errors. These shortcomings contribute to potentially error-prone and non-reproducible research.²

1
2
3 As the practice of medicine becomes more digitally-driven, because very few healthcare
4 provider organizations have data on large enough populations for even basic functions, such as
5 diagnosis and decision support, it will be increasingly in their interest to procure access to multi-
6 institutional data. Currently, when healthcare providers see short term gains through the sale of
7 data, they often cede this benefit to third parties with no data-driven knowledge accruing back.
8 At some point affordability of these multi-institutional datasets may become a challenging issue
9 for health care providers attempting to practice in a digital medicine ecosystem.³ If the goal of a
10 healthcare institution is simply commercial gain, the extant health data market might offer short
11 term advantages. If the goals are a learning health system, high quality research, and entry into
12 durable and transparent compacts with patients, there are better ways.
13
14
15
16
17
18
19
20
21
22
23
24
25

26 Broadly speaking, two approaches can address the torrential leak of de-identified health
27 record data. One is to establish best practice among the data providers—usually HIPAA covered
28 entities. The other is to strengthen legal and regulatory protections.
29
30
31
32

33 Healthcare institutions can strengthen protection of their interests and their patients'
34 privacy by treating de-identified data more like protected health information. First of all, patients
35 should be informed by their healthcare providers, through consents to treat and privacy notices,
36 that upon arrival to the hospital that their data may be used in a learning healthcare system, and
37 when appropriate shared with commercial parties.
38
39
40
41
42
43

44 Secondly, when data must be externalized, proper contractual controls should be
45 implemented to ensure that the data never passes outside the construct of the specific
46 arrangement, cannot be linked with other datasets, and that re-identification is prohibited. This
47 approach also enables disclosures to patients allowing full transparency of how the institution is
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 sharing and using their data. Ideally, the covered entity data provider is a genuine collaborator in
4
5 the effort, not just a data supplier.
6

7
8 Thirdly, a preferred approach to permitting use of data is so-called “behind the glass”
9
10 access where data do not leave the institution and instead the analytics are brought to the data.
11
12 One approach is to create shared analytic workspaces, or enclaves, where health systems can
13
14 interact directly with data and technology innovators but without the need to externalize patient
15
16 data. Data can be combined on a project-by-project basis, under protective contracts or data use
17
18 agreements. To aggregate multi-institutional datasets, federated analytic models permit data to
19
20 stay at the originator sites.
21
22

23
24 To improve legal and regulatory oversight, first of all, states, and the federal government,
25
26 should join California in making re-identification of de-identified health data illegal.⁴ This would
27
28 have an important effect on the market, but notably might not prevent a malicious actor from
29
30 exposing a patient’s medical information.
31
32

33
34 Secondly, it is worth examining closely the advantages and drawbacks of the European
35
36 Union’s Global Data Protection Regulations (GDPR) ‘Right to Erasure,’ which ensures a data
37
38 subject can choose to be erased from a dataset upon request, without undue delay when the data
39
40 are being used for purposes other than for which it was initially collected; the data subject
41
42 withdraws consent or objects to the use; the data are being used unlawfully, required for issues of
43
44 legal and/or regulatory compliance and other similar situations.⁵ Right to erasure shifts the
45
46 burden of risk from the patient back to institutions. However, if the patient does not opt out at the
47
48 beginning of the project, they should understand that their records cannot be subsequently
49
50 located in a truly de-identified dataset. In crafting laws or regulations inspired by GDPR
51
52
53
54
55
56
57
58
59
60

1
2
3 provisions, positive uses should be preserved. In a learning healthcare system for example, an
4
5 opt-out model might severely bias the dataset and prevent accurate analytics.
6

7
8 HIPAA and its privacy rule were crafted in the pre-electronic health record era. There is
9
10 now an opportunity for health systems, legislators, and regulators to protect health record data
11
12 beyond what HIPAA currently mandates, while promoting and supporting beneficial uses toward
13
14 improving health and optimizing health care delivery.
15
16
17
18

19 **References**

- 20
21 1. Rocher L, Hendrickx JM, de Montjoye Y-A. Estimating the success of re-identifications in
22
23 incomplete datasets using generative models. *Nat Commun* 2019;10(1):3069.
24
25
- 26
27 2. Mehra MR, Ruschitzka F, Patel AN. Retraction-Hydroxychloroquine or chloroquine with or
28
29 without a macrolide for treatment of COVID-19: a multinational registry analysis. *Lancet*
30
31 2020;395(10240):1820.
32
- 33
34 3. Mandl KD, Bourgeois FT. The Evolution of Patient Diagnosis: From Art to Digital Data-
35
36 Driven Science. *JAMA* 2017;318(19):1859–60.
37
- 38
39 4. California Legislature Adopts CCPA Exemption for Information Deidentified in
40
41 Accordance with the HIPAA Privacy Rule [Internet]. 2020 [cited 2021 Jan 19]; Available
42
43 from: [https://www.insideprivacy.com/ccpa/california-legislature-adopts-ccpa-exemption-](https://www.insideprivacy.com/ccpa/california-legislature-adopts-ccpa-exemption-for-information-deidentified-in-accordance-with-the-hipaa-privacy-rule/)
44
45 [for-information-deidentified-in-accordance-with-the-hipaa-privacy-rule/](https://www.insideprivacy.com/ccpa/california-legislature-adopts-ccpa-exemption-for-information-deidentified-in-accordance-with-the-hipaa-privacy-rule/)
46
- 47
48 5. Bovenberg J, Peloquin D, Bierer B, Barnes M, Knoppers BM. How to fix the GDPR's
49
50 frustration of global biomedical research. *Science* 2020;370(6512):40–2.
51
52
53
54
55
56
57
58
59
60