

Data

Day 2 — The foundation, constraint, and failure mode of health AI

TOM POLLARD · MIT · JUNE 23, 2026

Where we are on the spine

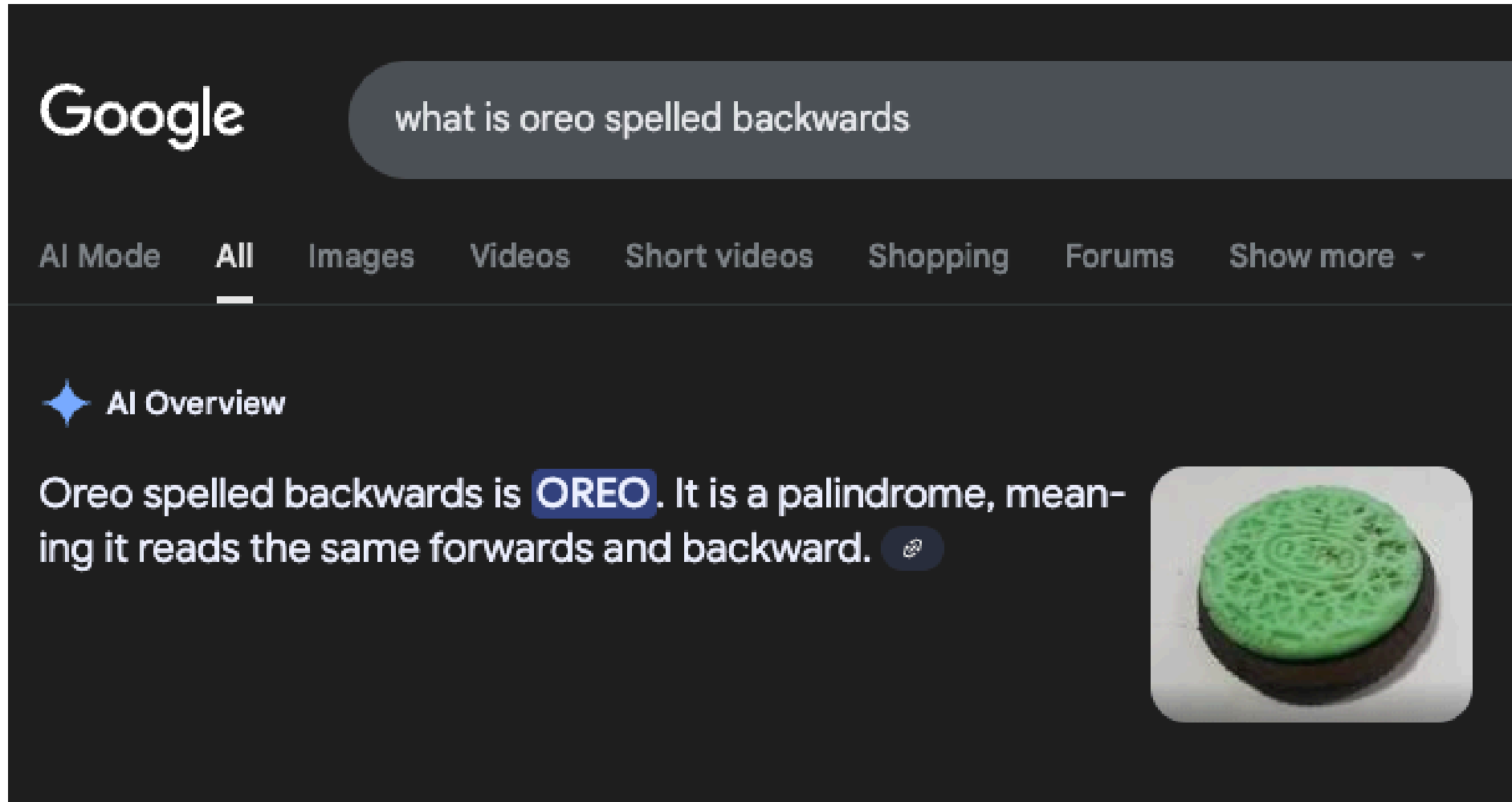
- **Yesterday:** an interesting topic → a well-posed problem.
- **Today:** a problem → the data that *actually exists* to attack it.

[↩ Recap problem statements](#)

The last decade was dominated by models.

The next will be shaped by data, governance, and deployment.

Are we okay with this?



The image shows a screenshot of a Google search interface. The search bar contains the text "what is oreo spelled backwards". Below the search bar, the "All" tab is selected. The search results include an "AI Overview" section with a blue star icon. The text in the AI Overview states: "Oreo spelled backwards is **OREO**. It is a palindrome, meaning it reads the same forwards and backward." To the right of this text is a small image of a single Oreo cookie with green frosting. At the bottom left of the screenshot, the text "LEARNING AND INFERENCE" is visible.

Why does Google think Oreo is a palindrome?

Google



why does google think oreo is a palindrome

✦ AI Overview

Google thinks "Oreo" is a palindrome **because of a famous flaw in how its AI processes text**. Rather than reading letter-by-letter, the AI reads words as chunks, leading to a few specific "glitch" reasons: [LinkedIn · Shanmugam Natarajan +2](#)

- **Tokenization Confusion:** AI models use "tokens" (chunks of text) instead of individual letters. When it scrambles or reverses these tokens, the logic breaks, causing the AI to confidently hallucinate that O-R-E-O reversed is still O-R-E-O. [LinkedIn · Shanmugam Natarajan +3](#)
- **Internet Sarcasm:** The claim became a viral internet meme. Because Google's AI generates answers based on top search results, it frequently "learns" the joke from social media and Reddit threads rather than from a dictionary. [LinkedIn · Gregor Berkowitz +4](#)
- **Visual/Structural Logic:** Some online jokes and AI models defend the idea structurally. Since an Oreo is essentially "cookie, cream, cookie," reversing the concept maintains the same pattern. [Facebook · A different type of Art +4](#)

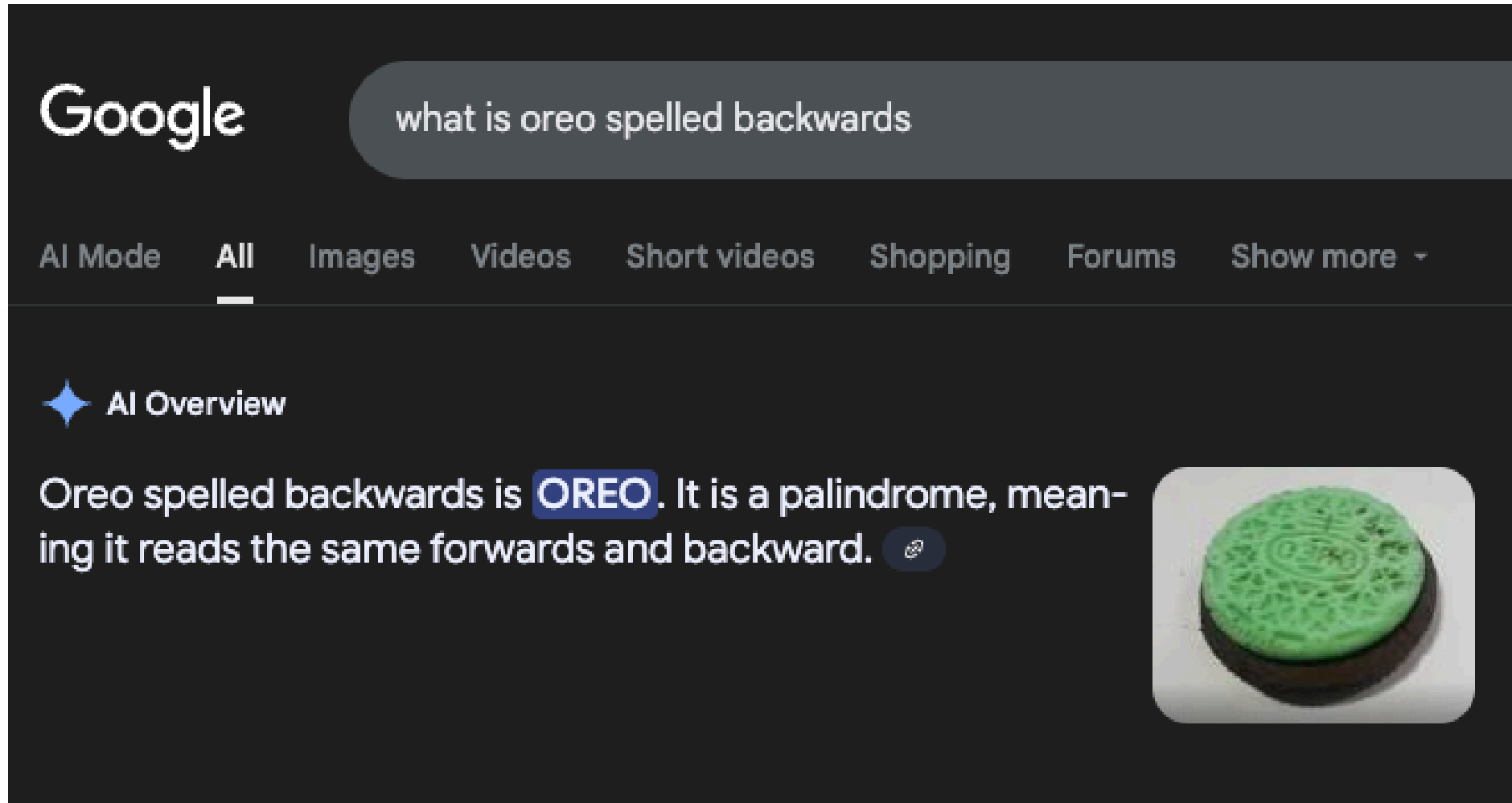
They fixed it!

The image shows a screenshot of a Google search interface. At the top left is the Google logo. To its right is a search bar containing the text "What is oreo spelled backwards?". Below the search bar are navigation tabs: "AI Mode", "All", "Images", "Videos", "Shopping", "Forums", and "Short videos". The "All" tab is currently selected. Below the tabs, there is a section titled "AI Overview" with a blue star icon. The text in this section reads: "Oreo spelled backwards is **oero**.  Reddit · r/notinteresting". Below this, there is another paragraph: "While it is a popular misconception that it spells the exact same word forwards and backwards, flipping the letters 'O-R-E-O' in reverse produces 'O-E-R-O'.  Reddit · r/words".

They fixed it?

The image is a screenshot of a Google search interface. At the top left is the Google logo. To its right is a search bar containing the text "What well known cookie name is a palindrome?". Below the search bar is a navigation bar with tabs for "AI Mode", "All", "Images", "Forums", "Shopping", "Videos", and "Short videos". The "All" tab is currently selected. Below the navigation bar, there is a section titled "AI Overview" with a blue star icon. The main text of the AI Overview reads: "The most well-known cookie name that is a palindrome (spelled the same forwards and backwards) is **Oreo**." To the right of this text is a small red icon and the text "Reddit · r/GoogleAIGoneWild".

Are we okay with this in healthcare?



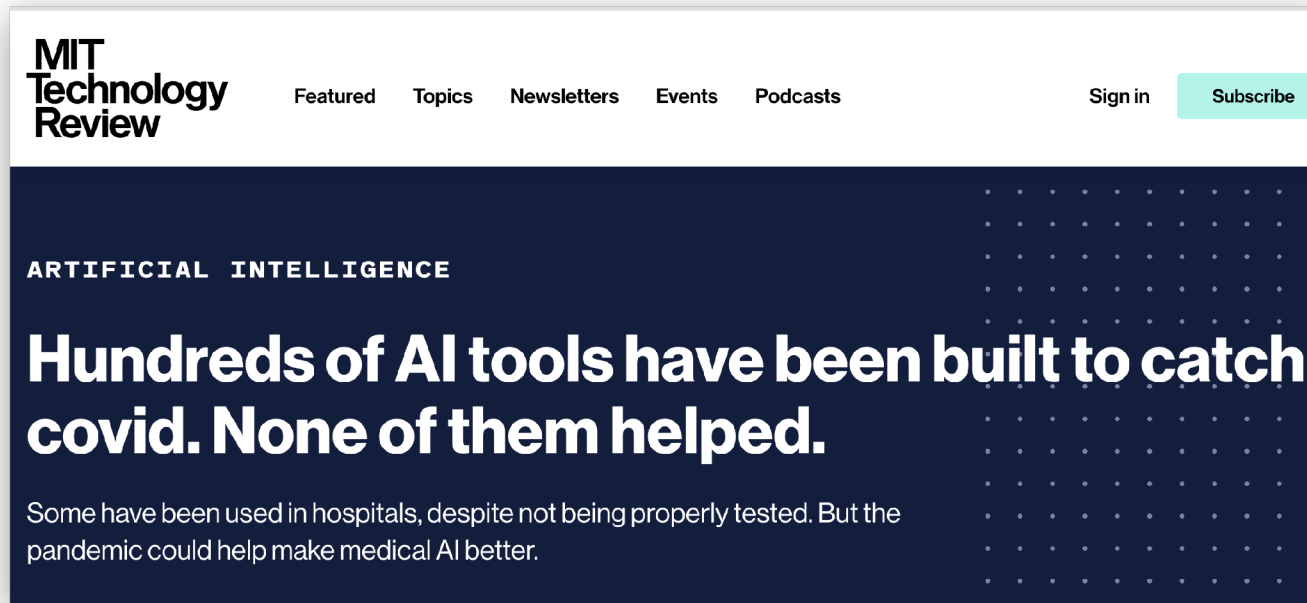
The image shows a screenshot of a Google search interface. The search bar contains the text "what is oreo spelled backwards". Below the search bar, the "All" tab is selected. The search results include an "AI Overview" section with a blue star icon. The text in the AI Overview states: "Oreo spelled backwards is **OREO**. It is a palindrome, meaning it reads the same forwards and backward." To the right of this text is a small image of a single Oreo cookie with green frosting. At the bottom left of the screenshot, the text "LEARNING AND INFERENCE" is visible.

Public health data has enabled reproducible research, benchmarks, education, and community science.

**Has it enabled
trustworthy,
meaningful, and
deployable AI for
health?**

COVID was a stress test

Researchers around the world turned their focus to building predictive models to beat covid.



The image is a screenshot of the MIT Technology Review website. At the top left is the MIT Technology Review logo. To its right are navigation links: 'Featured', 'Topics', 'Newsletters', 'Events', and 'Podcasts'. Further right are 'Sign in' and a teal 'Subscribe' button. Below the navigation is a dark blue banner with a grid of small white dots. The banner contains the text: 'ARTIFICIAL INTELLIGENCE' in all caps, followed by the headline 'Hundreds of AI tools have been built to catch covid. None of them helped.' in large white font. Below the headline is a sub-headline: 'Some have been used in hospitals, despite not being properly tested. But the pandemic could help make medical AI better.'

AHLI HEALTH AI SUMMER CAMP 2026

What went wrong?

"...**datasets spliced together** from multiple sources and [containing] duplicates."



What did the algorithms learn?

"were tested on the same data they were trained on...learnt to identify duplicate patients"

"learnt to use the text font to make predictions".

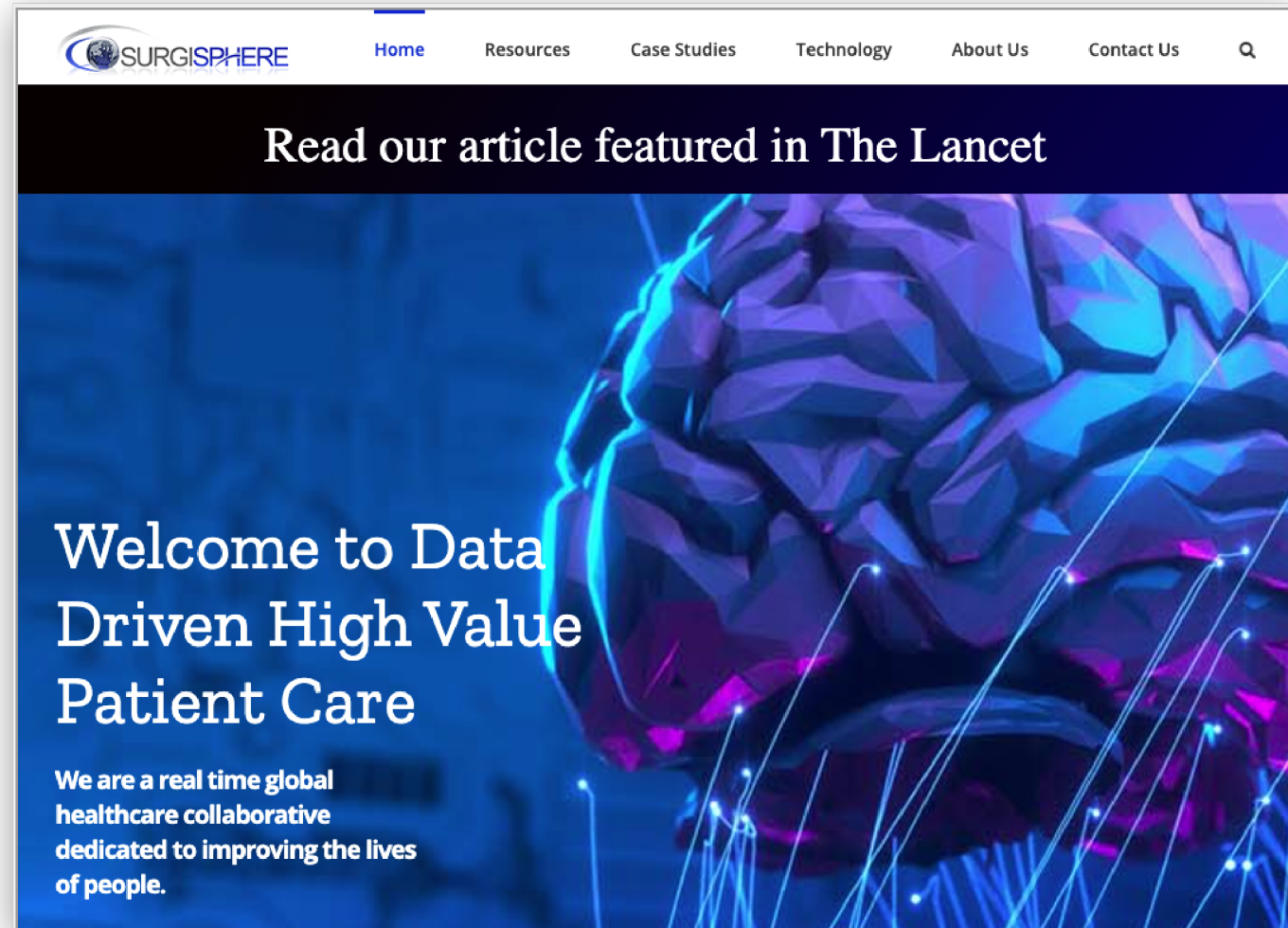
"learned to predict risk from a person's position"



Surgisphere

- Led by Sapan Desai, Vascular Surgeon (MD/PhD Chicago)
- 100,000 patients across 671 hospitals
- Papers in NEJM and The Lancet
- Huge global impact

Surgisphere website from archive.today (May 2020)



Provenance of data questioned

- Papers retracted!
- Huge waste of time
- Caused confusion and distrust

The New York Times

New Covid Shot Questions

Symptoms and Treatment

Free Tests Are Back

New Vaccines Are Coming

Two Huge Covid-19 Studies Are Retracted After Scientists Sound Alarms

The reports, published in two leading journals, were retracted after authors could not verify an enormous database of medical records.

FINANCIAL TIMES

Pharmaceuticals sector [+ Add to myFT](#)

WHO says it did not see Surgisphere data that halted virus drug trial

Authors retract Lancet and NEJM studies after company declines to transfer data sets





This ~~was~~ a turning point

This *should have been* a turning point

Evidence of unreliable data and poor data provenance in clinical prediction model research and clinical practice

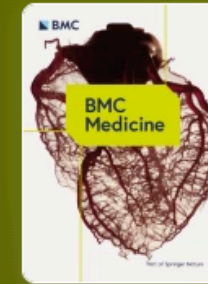
Research | [Open access](#) | Published: 04 June 2026

[Cite this article](#)

✓ You have full access to this [open access](#) article

[Download PDF](#) 

 [Save article](#)



[BMC Medicine](#)

[Aims and scope](#) →

[Submit manuscript](#) →

[Alexander D Gibson](#) , [Nicole M White](#), [Gary S Collins](#) & [Adrian G Barnett](#)

Abstract

Background

Clinical prediction models are often created using large routinely collected datasets. It is essential that prediction models are developed with appropriate data and methods and transparently reported to ensure that decisions are based on reliable predictions. Kaggle is a popular competition and data repository website where users learn and apply analysis skills on a range of datasets.

Methods

We identified two large, publicly available Kaggle datasets, on stroke and diabetes, that lack clear data provenance, but are widely used in clinical prediction models in peer reviewed

Sections

[Abstract](#)

[Abbreviations](#)

[Acknowledgements](#)

[Funding](#)

[Author information](#)

[Ethics declarations](#)

[Additional information](#)

[Supplementary Information](#)

[Rights and permissions](#)

[About this article](#)



These failures were not just “bad data”

They are infrastructural:

- provenance — where did the data come from?
- cohort logic — who was included, and why?
- leakage control — what information was available when?
- evaluation — what artifacts did the model learn?
- governance — who was allowed to use the data, and how?

No one interrogated how the data came to exist and whether it was the right data for the problem.

AHLI HEALTH AI SUMMER CAMP 2026
THE ORIGIN OF THE DATA

The data-generating process

Byproduct of care vs. designed for research

Byproduct of care

- Generated while clinicians treat, document, and bill.
- Abundant, messy, biased by workflow.
- Most clinical data lives here.

Designed for research

- Generated by a protocol to answer a defined question.
- Rarer, cleaner *for its purpose* — still biased.
- Benchmarks, registries, trials.

We need to ask:

- what process created the data, and
- what did that process make visible or invisible?

Data is **not** a neutral input.

It is generated by a process, for a purpose — and that shapes what you can credibly ask.

Patient timeline

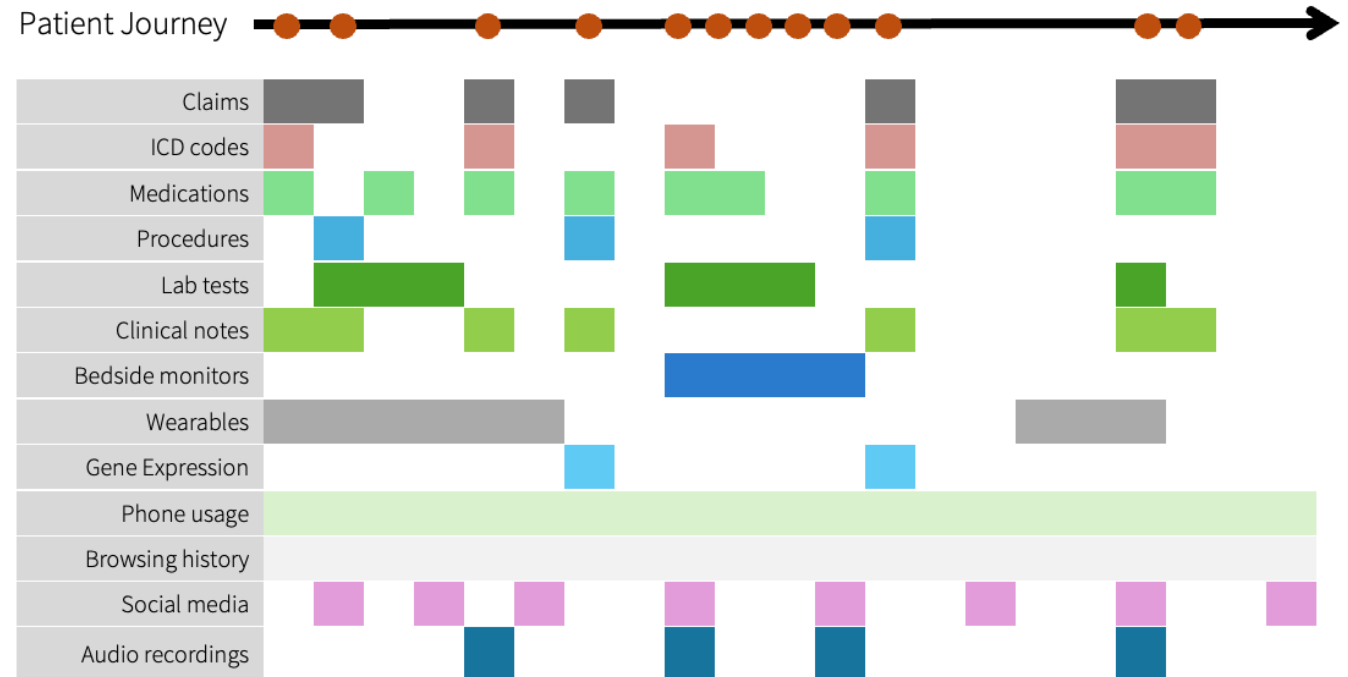
Often we think about isolated sets of data. ICU data; OR data; wearables data.

Helpful to remind ourselves that we are dealing with a sequence of events

Times of health and sickness.

Typically we are viewing just snapshots of the journey.

Credit: Emily Alsentzer.



Case study — how a clinical dataset comes to exist

MIMIC as an example:

- Deidentified critical-care data from patients admitted to a Boston Hospital
- **Built** through extraction, structuring, and a substantial de-identification effort.
- A "ready-to-use" dataset represents **significant invisible work** — and the choices in that work shape the data.



Critical care



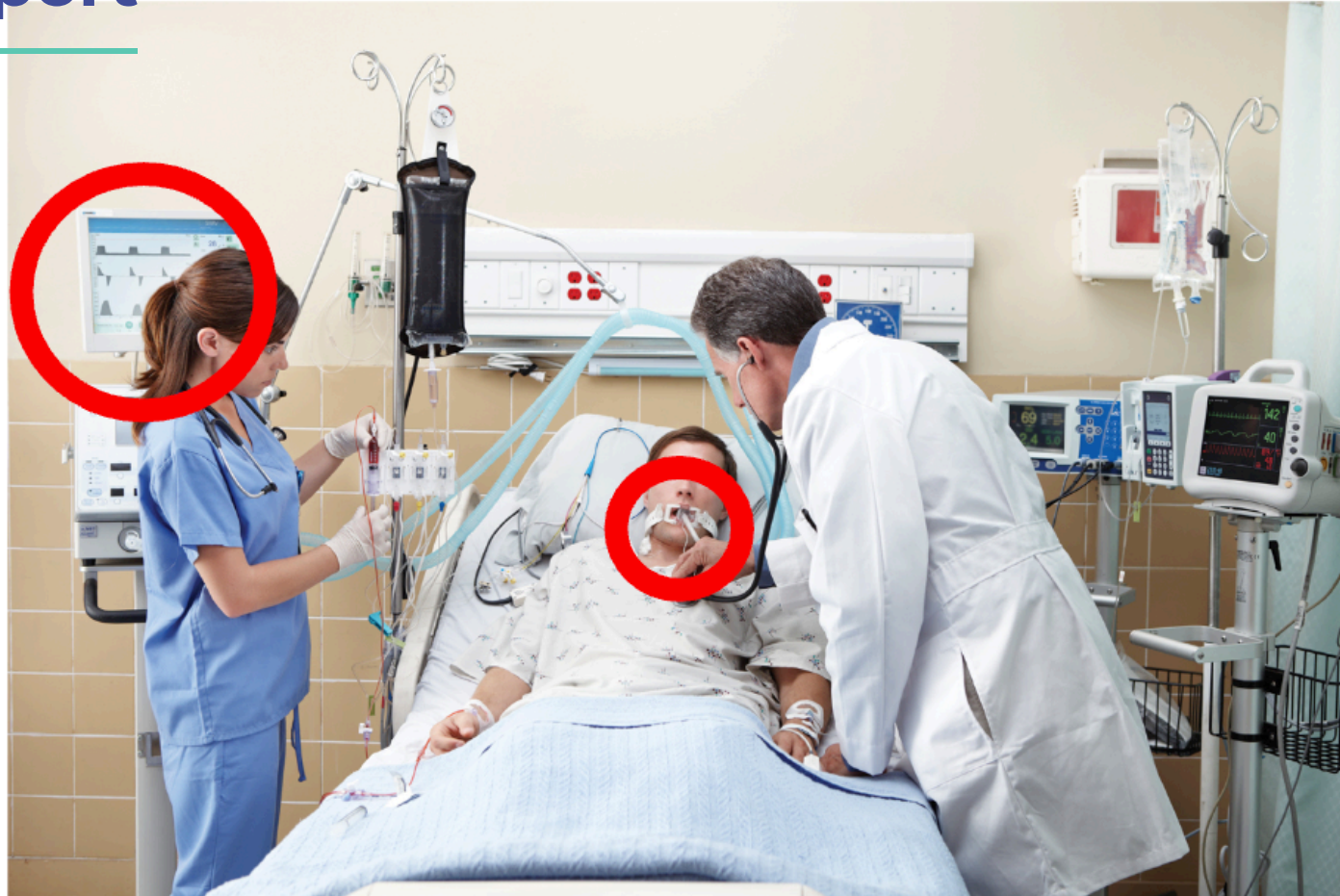
Infusions



Intermittent treatment



Organ support



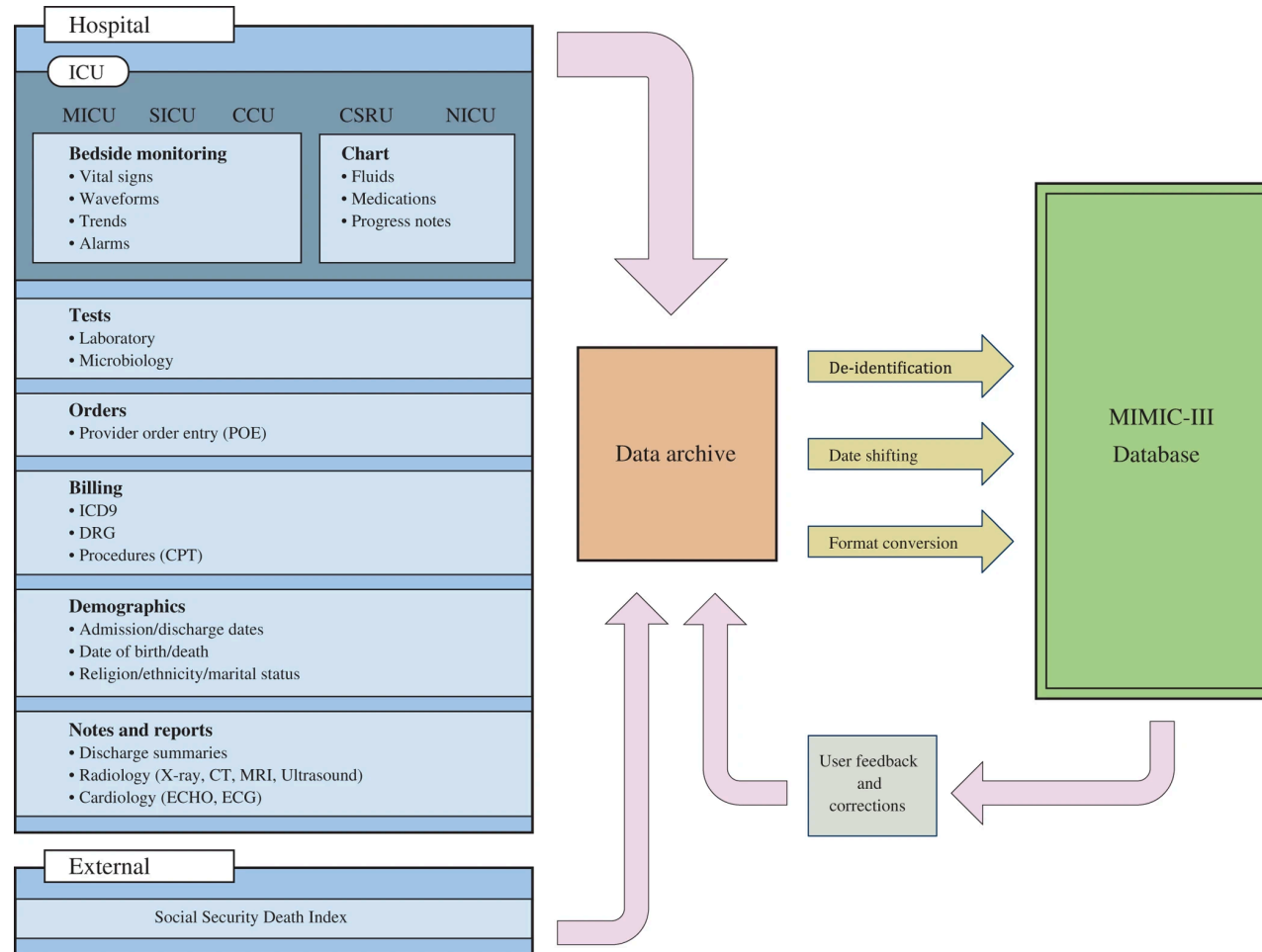
Waveforms



Observations, discussion, documentation



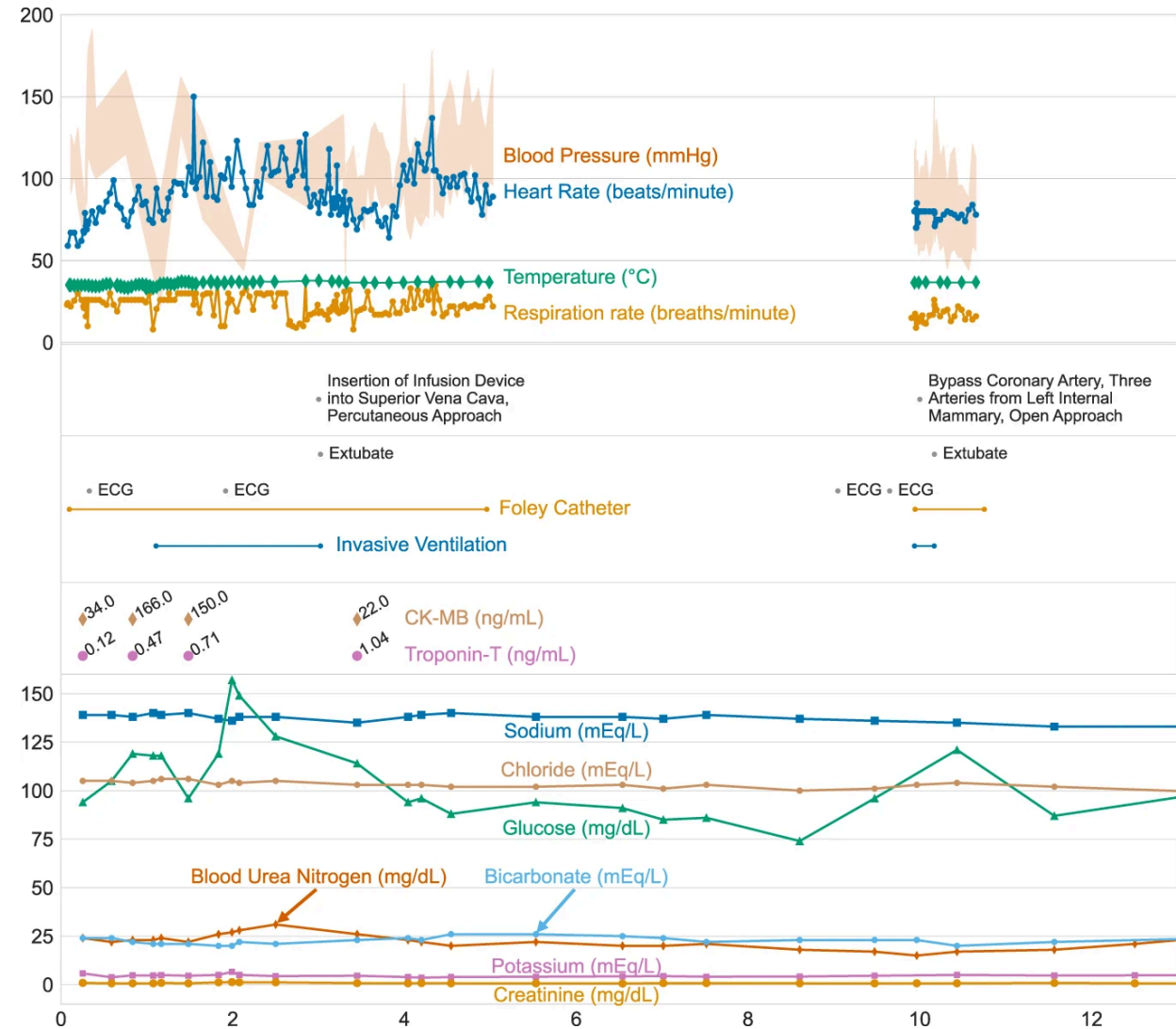
MIMIC



hadm_id=28503629

Vital signs
captured only
during ICU stay

Target monitoring
increases
frequency of
measurements



Real-life data



r/MachineLearning

Find a community, post, or user

LOG IN

↑ [captkrob](#) 2 points · 1 year ago

↓ If you want an idea of a "real world" disease dataset, try looking at the MIMIC-III database

(<https://mimic.physionet.org/>).

It's probably significantly more work to request the data and set it up in an easy format for ML than you had ever planned on, but this is the kind of thing actual data scientists and informaticians working on this kind of problem are used to dealing with. The dataset is tremendously rich, but also incredibly noisy. In other words, much like medical practice.

Share Save

Cleaning should be visible and reusable

Encourage data processing pipelines, cohort selection to be shared in a public code repository.

Research and Applications

The MIMIC Code Repository: enabling reproducibility in critical care research

Alistair EW Johnson,^{1,*} David J Stone,² Leo A Celi,^{1,3} and Tom J Pollard¹

¹Massachusetts Institute of Technology, Cambridge, MA, USA, ²University of Virginia School of Medicine, Charlottesville, VA, USA, and ³Beth Israel Deaconess Medical Center, Boston, MA, USA

*Corresponding Author: Alistair EW Johnson, Massachusetts Institute of Technology, E25-505, 77 Massachusetts Avenue, Cambridge, MA 02139, USA. E-mail: aewj@mit.edu.

Received 25 May 2017; Revised 11 July 2017; Editorial Decision 21 July 2017; Accepted 27 July 2017

ABSTRACT

Objective: Lack of reproducibility in medical studies is a barrier to the generation of a robust knowledge base to support clinical decision-making. In this paper we outline the Medical Information Mart for Intensive Care (MIMIC) Code Repository, a centralized code base for generating reproducible studies on an openly available critical care dataset.

Materials and Methods: Code is provided to load the data into a relational structure, create extractions of the data, and reproduce entire analysis plans including research studies.

Care data are selected by policy

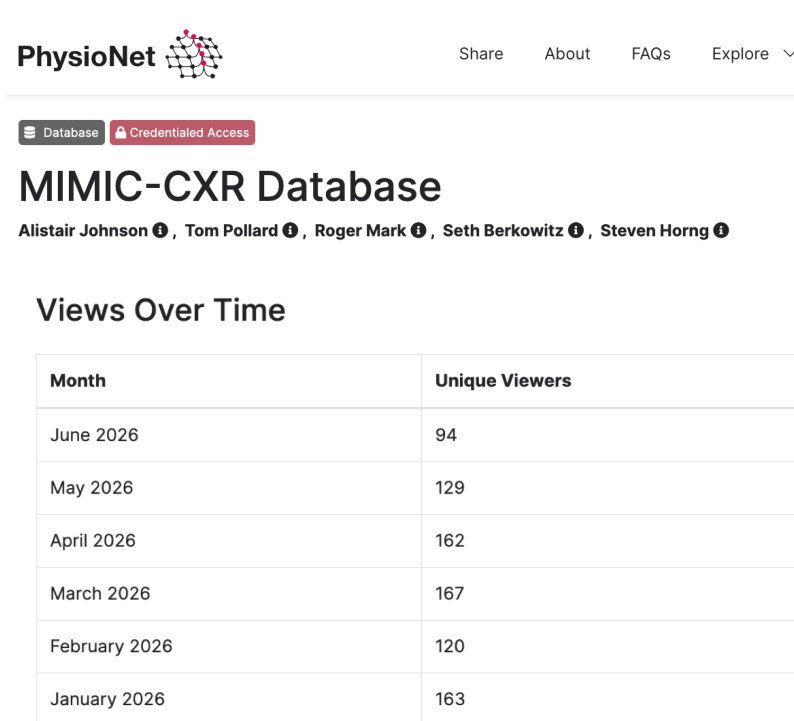
- Missingness is **MNAR by construction** — a value is absent because of a *decision* (not ordered ⇒ not worried).
- In EHR data, “controls” often mean patients who were not tested, treated, coded, or escalated, not a random sample of true negatives.
- Conditioning on a measurement (e.g. "lab was drawn") leads to **sampling bias**.

The data encodes the *policy* that produced it — and policies are not exchangeable across sites or time.

Each modality captures measurement *decisions*

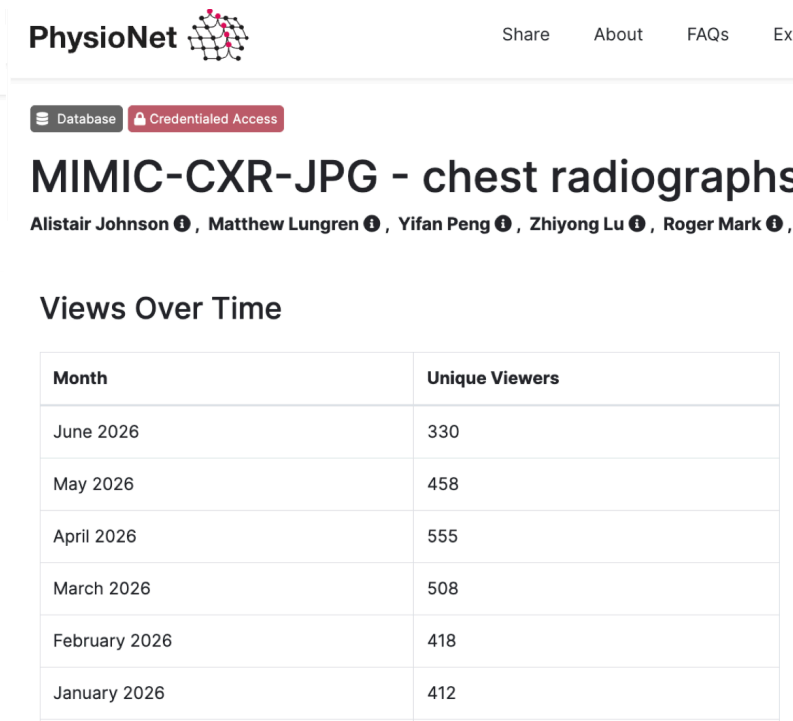
MODALITY	THE MEASUREMENT IS ITSELF A DECISION...	...WHICH BREAKS
Labs	ordered when a clinician is concerned	condition on "drawn" → sampling bias
Vitals	sampled faster when acuity is high	frequency is a severity proxy, not noise
Notes	dictated selectively, templated	presence ≠ incidence; documentation drift
Dx / billing codes	coded for reimbursement	prevalence tracks coding incentives
Imaging	protocol- and referral-gated	case-mix bias; site-specific acquisition

Design choices influence downstream use



The screenshot shows the PhysioNet website for the MIMIC-CXR Database. The page includes a navigation bar with 'Share', 'About', 'FAQs', and 'Explore'. Below the navigation, there are buttons for 'Database' and 'Credentialed Access'. The main title is 'MIMIC-CXR Database' with authors listed below it. A 'Views Over Time' section contains a table with the following data:

Month	Unique Viewers
June 2026	94
May 2026	129
April 2026	162
March 2026	167
February 2026	120
January 2026	163



The screenshot shows the PhysioNet website for the MIMIC-CXR-JPG - chest radiographs dataset. The page includes a navigation bar with 'Share', 'About', 'FAQs', and 'Exp'. Below the navigation, there are buttons for 'Database' and 'Credentialed Access'. The main title is 'MIMIC-CXR-JPG - chest radiographs' with authors listed below it. A 'Views Over Time' section contains a table with the following data:

Month	Unique Viewers
June 2026	330
May 2026	458
April 2026	555
March 2026	508
February 2026	418
January 2026	412

AHLI HEALTH AI SUMMER CAMP 2026
GETTING THE DATA

De-identification, governance, and access

De-identification and governance

- Data is shared under Health Insurance Portability and Accountability Act of 1996 (Safe Harbor's 18 identifiers, expert determination).

"A major goal of the Privacy Rule is to assure that individuals' health information is properly protected while allowing the flow of health information needed to provide and promote high quality health care and to protect the public's health and well being." .. "The Rule strikes a balance that permits important uses of information, while protecting the privacy of people who seek care and healing."

"There are no restrictions on the use or disclosure of de-identified health information. De-identified health information neither identifies nor provides a reasonable basis to identify an individual."

1. Manual review (time-consuming)

- Scan infrequent values in string fields (e.g. group lab tests, review tail)
- Scrape news articles



Injured woman hit by ball released from hospital



By Alec Shirkey
July 11, 2015



BOSTON – Fenway Park became the scene of another scary situation during the Red Sox's **5-1 loss** to the Yankees on Friday night when a woman was injured by a foul ball.

A Beth Israel Deaconess spokesman said the hurt woman was treated and released from the hospital on Saturday.

2. Pattern matching (predictable, fragile)

- Known identifiers, e.g. match variations on:
 - Patient name and location
 - Dates
 - "VIP list"
- Document structure:
 - Fields like "Name: XXX"
- PHI structure:
 - Dates (YYYY/MM/DD)
 - Email (username@domain.tld)

3. Machine learning (unpredictable, flexible)

Deidentification of free-text medical records using pre-trained bidirectional transformers

Alistair E. W. Johnson

aewj@mit.edu

Massachusetts Institute of Technology
Cambridge, MA, USA

Lucas Bulgarelli

lucas1@mit.edu

Massachusetts Institute of Technology
Cambridge, MA, USA

Tom J. Pollard

tpollard@mit.edu

Massachusetts Institute of Technology
Cambridge, MA, USA

ABSTRACT

The ability of caregivers and investigators to share patient data is fundamental to many areas of clinical practice and biomedical research. Prior to sharing, it is often necessary to remove identifiers such as names, contact details, and dates in order to protect patient privacy. Deidentification, the process of removing identifiers, is challenging, however. High-quality annotated data for developing models is scarce; many target identifiers are highly heterogeneous (for example, there are uncountable variations of patient names); and in practice anything less than perfect sensitivity may be considered a failure. As a result, patient data is often withheld when sharing would be beneficial, and identifiable patient data is often

1 INTRODUCTION

The advent of large, open access text corpuses and the resurgence of neural networks has driven advances in state-of-the-art model performance in natural language processing [10, 27]. Barriers to sharing clinical text, however, have stifled progress in the medical domain. An unintended consequence is that research has become hyperfocused on the few datasets that are readily accessible. MIMIC-III, one of the only public sources of electronic health record data, for example, has been referred to as “one of the most (over)analyzed clinical datasets” [12, 29]. This paucity of data is to the detriment of important issues including bias, generalizability, and reproducibility [5].

Deidentification is not bias-free

In the Name of Fairness: Assessing the Bias in Clinical Record De-identification

Yuxin Xiao*
yuxin102@mit.edu
Massachusetts Institute of Technology
USA

Tom Joseph Pollard
tpollard@mit.edu
Massachusetts Institute of Technology
USA

Shulammit Lim*
shulim@mit.edu
Massachusetts Institute of Technology
USA

Marzyeh Ghassemi
mghassem@mit.edu
Massachusetts Institute of Technology
USA

ABSTRACT

Data sharing is crucial for open science and reproducible research, but the legal sharing of clinical data requires the removal of protected health information from electronic health records. This process, known as de-identification, is often achieved through the use of machine learning algorithms by many commercial and open-source systems. While these systems have shown compelling results on average, the variation in their performance across different demographic groups has not been thoroughly examined. In this work, we investigate the bias of de-identification systems on names in clinical notes via a large-scale empirical analysis. To achieve this, we create 16 name sets that vary along four demographic dimensions: gender, race, name popularity, and the decade of popularity. We insert these names into 100 manually curated clinical

ACM Reference Format:

Yuxin Xiao, Shulammit Lim, Tom Joseph Pollard, and Marzyeh Ghassemi. 2023. In the Name of Fairness: Assessing the Bias in Clinical Record De-identification. In *2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT '23)*, June 12–15, 2023, Chicago, IL, USA. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3593013.3593982>

1 INTRODUCTION

The increased availability of clinical datasets [51, 73, 74] plays a significant role in the recent advancements in machine learning (ML)-aided healthcare systems [15, 39, 107, 118]. In order to share clinical trial data legally, stakeholders must adhere to the Health Insurance Portability and Accountability Act (HIPAA) Safe Harbor provisions by masking 18 types of protected health information (PHI). If done

Access mechanics

- 1 **Credentialing:** human-subjects training, approved identity.
- 2 **Data use agreement:** sign and comply with the DUA.

Realistic timelines: weeks to months.

Start early — before you need the data.

PhysioNet Share About FAQs Explore

Find, share, and reuse health data

[Publish Your Data](#) [Discover Datasets](#)

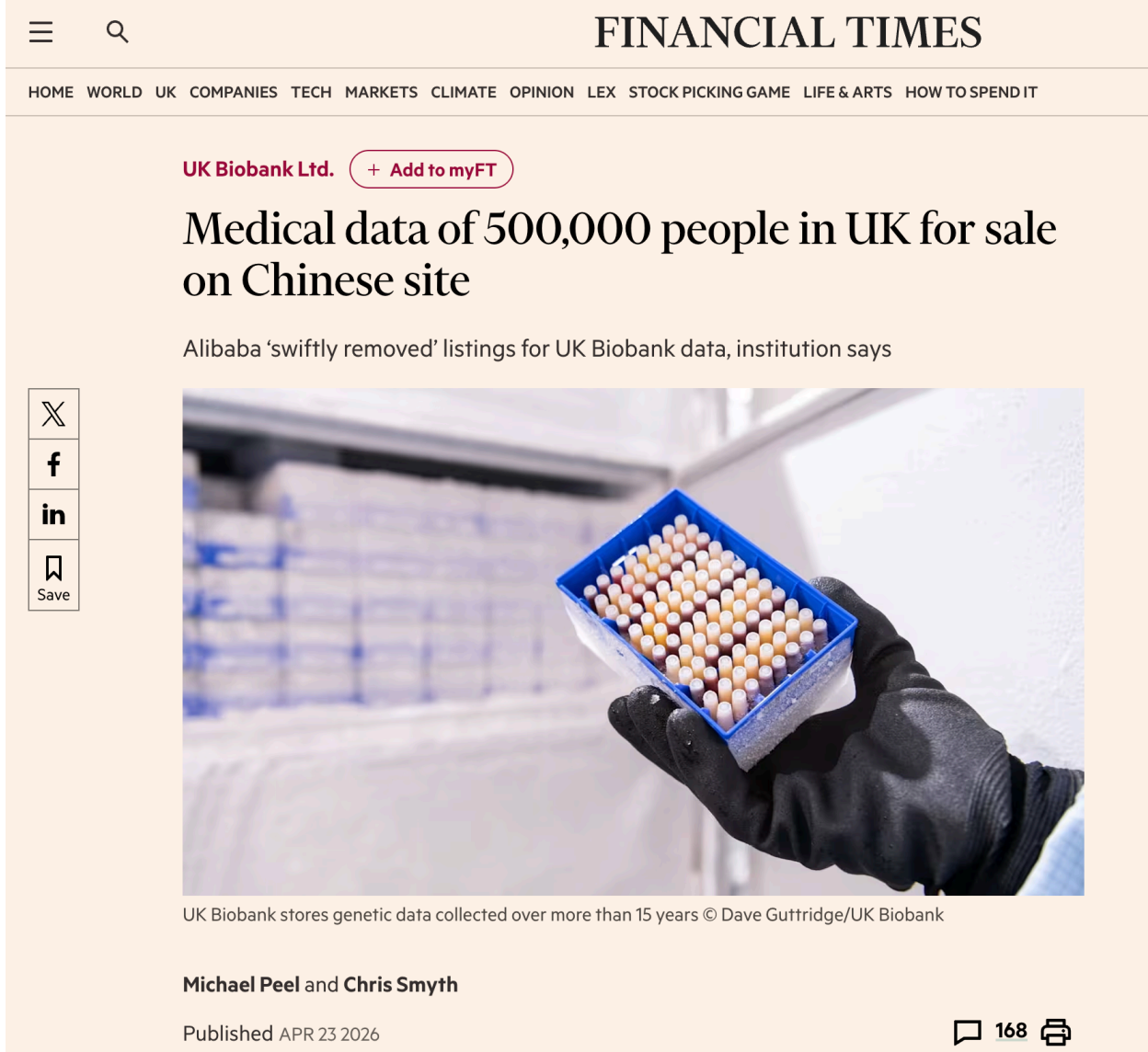
July 15, 2025
Access Restrictions Under DOJ Data Security Program
PhysioNet has introduced updated access policies for certain datasets to comply with the U.S. Department of Justice's Data Security Program (DSP) under Executive Order 14117. The DSP final rule took effect on April 8, 2025 and full enforcement began July 8, 2025: [https://www.justice.gov/opa/media/1396351/...](https://www.justice.gov/opa/media/1396351/)

Sept. 24, 2025
Use of MIMIC Data with Large Language Models and Online Services
We have received inquiries about the use of credentialed and restricted data on PhysioNet, including MIMIC-III, MIMIC-IV, MIMIC-CXR, and their derivatives, with large language models (LLMs) and online services. The PhysioNet Credentialed Data Use Agreement explicitly prohibits sharing access...

Aug. 18, 2025
Roger Mark and George Moody Receive the 2026 IEEE Biomedical Engineering Award
Each year, the IEEE Awards Board selects a distinguished group of individuals to receive IEEE's highest honors, recognizing exceptional achievements and significant contributions to technology, society, and the engineering profession.
We are honored to share that Professor Roge...

Increasing value of health data makes it a target

Even well governed resources face downstream misuse.



The screenshot shows the Financial Times website interface. At the top, there is a navigation bar with the 'FINANCIAL TIMES' logo and a search icon. Below it, a secondary navigation bar lists various sections: HOME, WORLD, UK, COMPANIES, TECH, MARKETS, CLIMATE, OPINION, LEX, STOCK PICKING GAME, LIFE & ARTS, and HOW TO SPEND IT. The main article is titled 'UK Biobank Ltd.' with a '+ Add to myFT' button. The headline reads 'Medical data of 500,000 people in UK for sale on Chinese site'. A sub-headline states 'Alibaba 'swiftly removed' listings for UK Biobank data, institution says'. To the left of the article is a vertical sidebar with social media sharing icons for Facebook, LinkedIn, and a 'Save' button. The main image shows a person wearing a black glove holding a blue tray filled with small, white, cylindrical vials. Below the image is a caption: 'UK Biobank stores genetic data collected over more than 15 years © Dave Guttridge/UK Biobank'. The authors are listed as 'Michael Peel and Chris Smyth'. At the bottom, it says 'Published APR 23 2026'. In the bottom right corner, there are icons for comments (168) and a print icon.

LLMS

Researchers want to send data to external services (e.g. LLMs) for their research.

This involves sharing with a 3rd party (not allowed under our Data Use Agreement)

Sept. 24, 2025

Use of MIMIC Data with Large Language Models and Online Services

We have received inquiries about the use of credentialed and restricted data on PhysioNet, including MIMIC-III, MIMIC-IV, MIMIC-CXR, and their derivatives, with large language models (LLMs) and online services. The PhysioNet Credentialed Data Use Agreement explicitly prohibits sharing access to the data with third parties, including sending it through APIs or using it on...



Membership attacks

High capacity models can memorize data. Membership attacks become a real threat.

External data: insurance claim: age, weekend accident, injuries

Released clinical data: de-identified admission: same age, same timing, same injuries

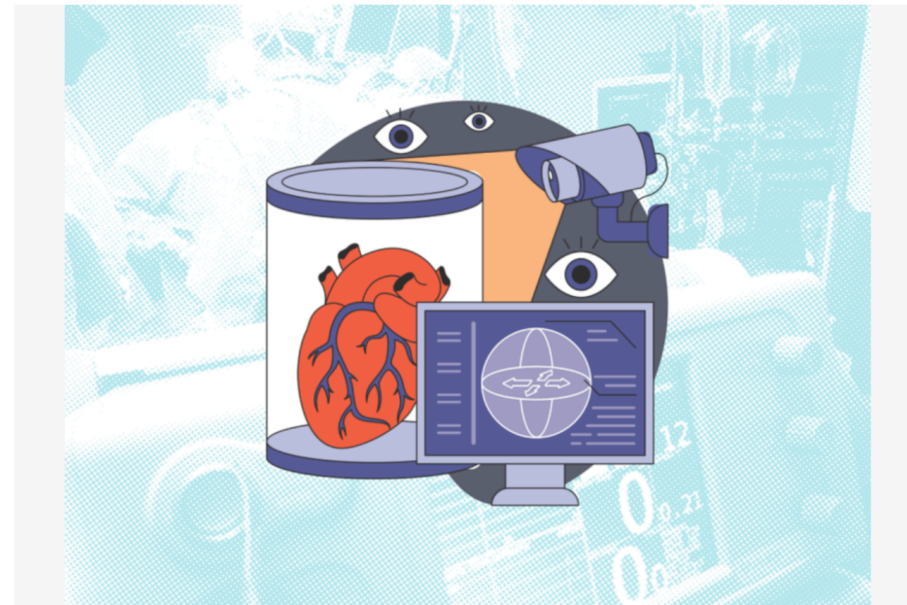
New information: clinical note shows suspected intoxication

MIT scientists investigate memorization risk in the age of clinical AI

New research demonstrates how AI models can be tested to ensure they don't cause harm by revealing anonymized patient health data.

Alex Ouyang | Abdul Latif Jameel Clinic for Machine Learning in Health
January 5, 2026

▼ PRESS INQUIRIES



MIT scientists are developing tests to ensure that AI models aren't memorizing sensitive patient information.

Image: Alex Ouyang/MIT Jameel Clinic, with Adobe Stock

Regulation is changing

Department of Justice recently introduced the Data Security Program

Limits sharing of deidentified health data with "Countries of Concern"

Preventing Access to U.S. Sensitive Personal Data and Government-Related Data by Countries of Concern or Covered Persons

A Rule by the [Justice Department](#) on 01/08/2025



PUBLISHED DOCUMENT: 2024-31486 (90 FR 1636)

- PDF
- Document Details
- Document Dates
- Table of Contents
- Related Documents
- Public Comments
- Regulations.gov Data
- Sharing

AGENCY:
National Security Division, Department of Justice.

SUMMARY:
The Department of Justice is issuing a final rule to implement [Executive Order 14117](#) of February 28, 2024 (Preventing Access to Americans' Bulk Sensitive Personal Data and United States Government-Related Data by Countries of Concern), by prohibiting and restricting certain data transactions with certain countries or persons.

DATES:
This rule has been classified as meeting the criteria under [5 U.S.C. 804\(2\)](#) and is effective April 8, 2025. However, at the conclusion of the Congressional review, if the effective date has been changed, the Department of Justice will publish a document in the **Federal Register** to establish the actual date of effectiveness or to terminate the rule. The incorporation by reference of certain material listed in this rule is approved by the Director of the Federal Register as of April 8, 2025.

Public data enables (weak?) research

Public data enables (weak?) research

Study explored trends in publications using public datasets

“excess of 12,000 papers in 2025”

“measures to control access to open data are required”



Quantifying new threats to health and biomedical literature integrity from rapidly scaled publications and problematic research

Matt Spick ^a, Anthony Onoja ^a, Charlie Harrison ^b, Stefan Stender ^c, Jennifer Byrne ^{d e}, Nophar Geifman ^a

Highlights

Key findings

- We estimate an annual excess of 12 thousand papers exploiting Open Science in 2025

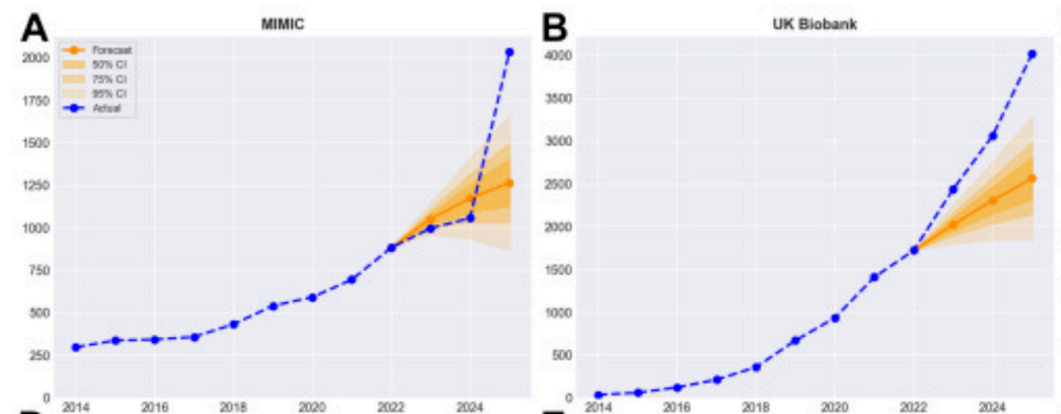
What this adds to what is known?

- Growth in AI-assisted low-quality paper production using Open Data has exploded - AI papers are introducing false discoveries and misleading claims to the literature

What is the implication and what should change now?

- Measures to deter formulaic submissions and control access to Open Data are required

Forecast vs Actual Trends (Selected Datasets)



AHLI HEALTH AI SUMMER CAMP 2026

Public datasets can create oversaturation

Healthcare, Machine Learning

The Shaky Foundations of Foundation Models in Healthcare

Scholars detail the current state of large language models in healthcare and advocate for better evaluation frameworks.

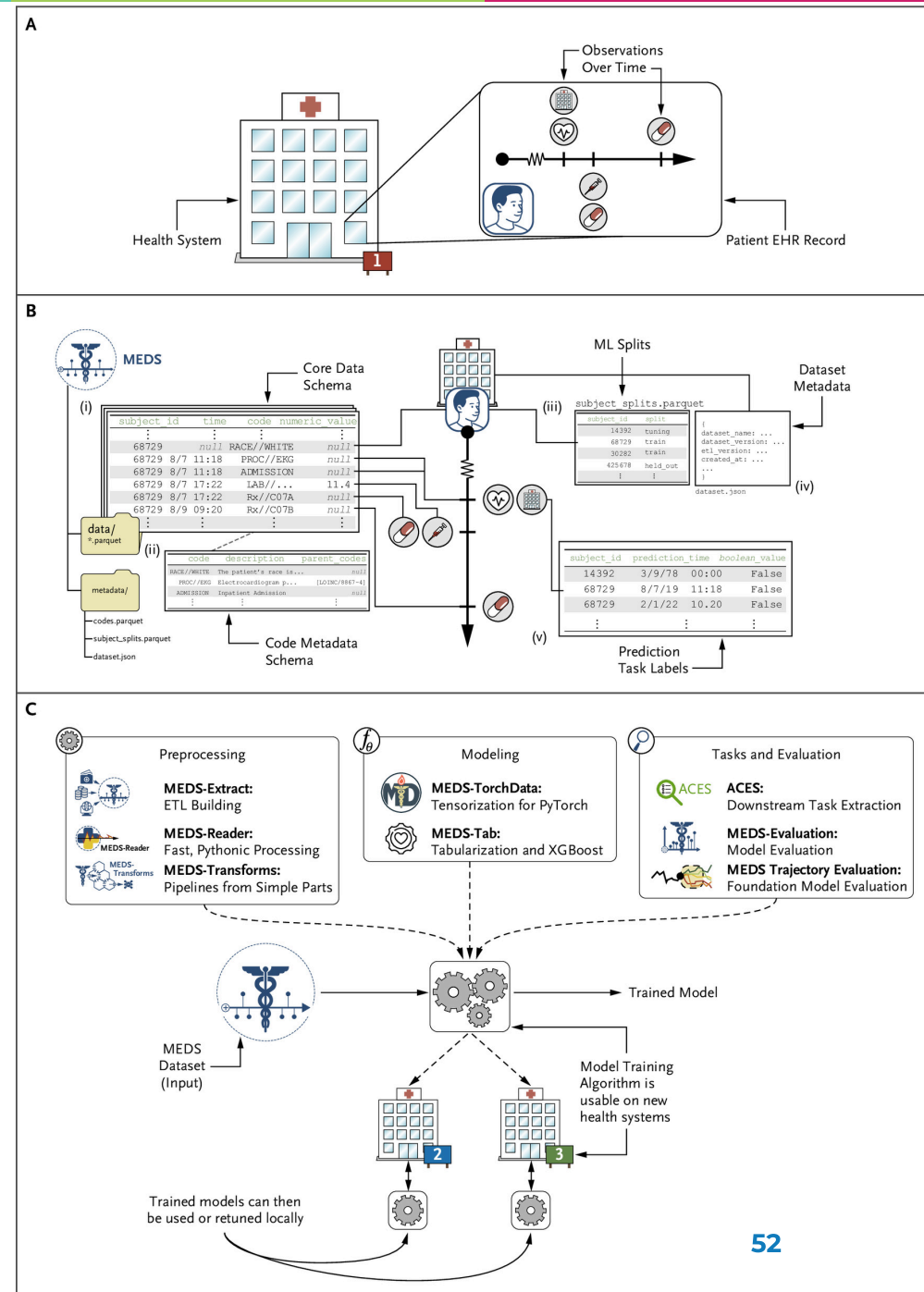
Feb 27, 2023 |

Michael Wornow, Yizhe Xu, Birju Patel, Rahul Thapa, Ethan Steinberg, Scott Fleming, Jason Fries, Nigam Shah



Standards solve syntax; representation is still open

- A common schema — **MEDS**, events as (patient, time, code, value) — makes data *portable*.
- **How to represent** irregular, multimodal, multi-scale event streams is an **open research problem** (e.g. tokenization, time encoding, and concept vocabularies).
- Standardization reduces busywork, allowing you to focus on the modeling question (Day 4).



AHLI HEALTH AI SUMMER CAMP 2026
READING BETWEEN THE LINES

What the data hides

Provenance-driven bias

- **Time:** coding, testing, and treatment change; the model learns the era.
- **Cohort definitions:** cases and controls encode how patients were found.
- **Selection:** the dataset contains who entered the system, not everyone at risk.
- **Documentation:** codes and notes reflect workflow, incentives, and physiology.

Ask: is this signal clinical, or is it a trace of how the data were produced?

When a label is a billing artifact

A diagnosis code looks like ground truth. It isn't.

- A code may appear because it **justifies reimbursement**, not because it's the primary clinical concern.
- It may be **absent** for a real condition that wasn't billed.
- Prevalence in the data tracks **coding incentives** as much as biology.





Train on that label, and you partly learn the hospital's billing rules.

↩ Day 1: label noise

→ Day 3: estimands

Measurements are not neutral

This Issue Views **25,689** | Citations **59** | Altmetric **2398**

 [Download PDF](#)  [More ▾](#)  [Cite This](#)  [Permissions](#)

Original Investigation

FREE

July 11, 2022

Assessment of Racial and Ethnic Differences in Oxygen Supplementation Among Patients in the Intensive Care Unit

Eric Raphael Gottlieb, MD, MS^{1,2,3}; Jennifer Ziegler, MD, MSc⁴; Katharine Morley, MD, MPH^{2,5}; [et al](#)

[» Author Affiliations](#) | [Article Information](#)

JAMA Intern Med. 2022;182(8):849-858. doi:10.1001/jamainternmed.2022.2587

 [Editorial Comment](#)

 [Related Articles](#)

Key Points

Question Are there differences in supplemental oxygen administration among patients of different races and ethnicities associated with pulse oximeter performance discrepancies?

Findings In this cohort study of 3069 patients in the intensive care unit, Asian, Black, and

AHLI HEALTH AI SUMMER CAMP 2026

X-rays encode race





THE LANCET Digital Health

[This journal](#) [Journals](#) [Publish](#) [Clinical](#) [Global health](#) [Multimedia](#) [Events](#) [About](#)

ARTICLES | [VOLUME 4, ISSUE 6, E406-E414, JUNE 2022](#)

[Download Full Issue](#)

AI recognition of patient race in medical imaging: a modelling study

[Judy Wawira Gichoya, MD](#)   • [Imon Banerjee, PhD](#) • [Ananth Reddy Bhimoreddy, MS](#) • [John L Burns, MS](#) • [Leo Anthony Celi, MD](#) • [Li-Ching Chen, BS](#) • [Ramon Correa, BS](#) • [Natalie Dullerud, MS](#) • [Marzyeh Ghassemi, PhD](#) • [Shih-Cheng Huang](#) • [Po-Chih Kuo, PhD](#) • [Matthew P Lungren, MD](#) • [Lyle J Palmer, PhD](#) • [Brandon J Price, MD](#) • [Saptarshi Purkayastha, PhD](#) • [Ayis T Pyrros, MD](#) • [Lauren Oakden-Rayner, MD](#) • [Chima Okechukwu, MS](#) • [Laleh Seyyed-Kalantari, PhD](#) • [Hari Trivedi, MD](#) • [Ryan Wang, BS](#) • [Zachary Zaiman](#) • [Haoran Zhang, MS](#) • [Show less](#)

Open Access • Published: May 11, 2022 • DOI: [https://doi.org/10.1016/S2589-7500\(22\)00063-2](https://doi.org/10.1016/S2589-7500(22)00063-2)

 Check for updates

Summary

Introduction

Methods

Results

Discussion

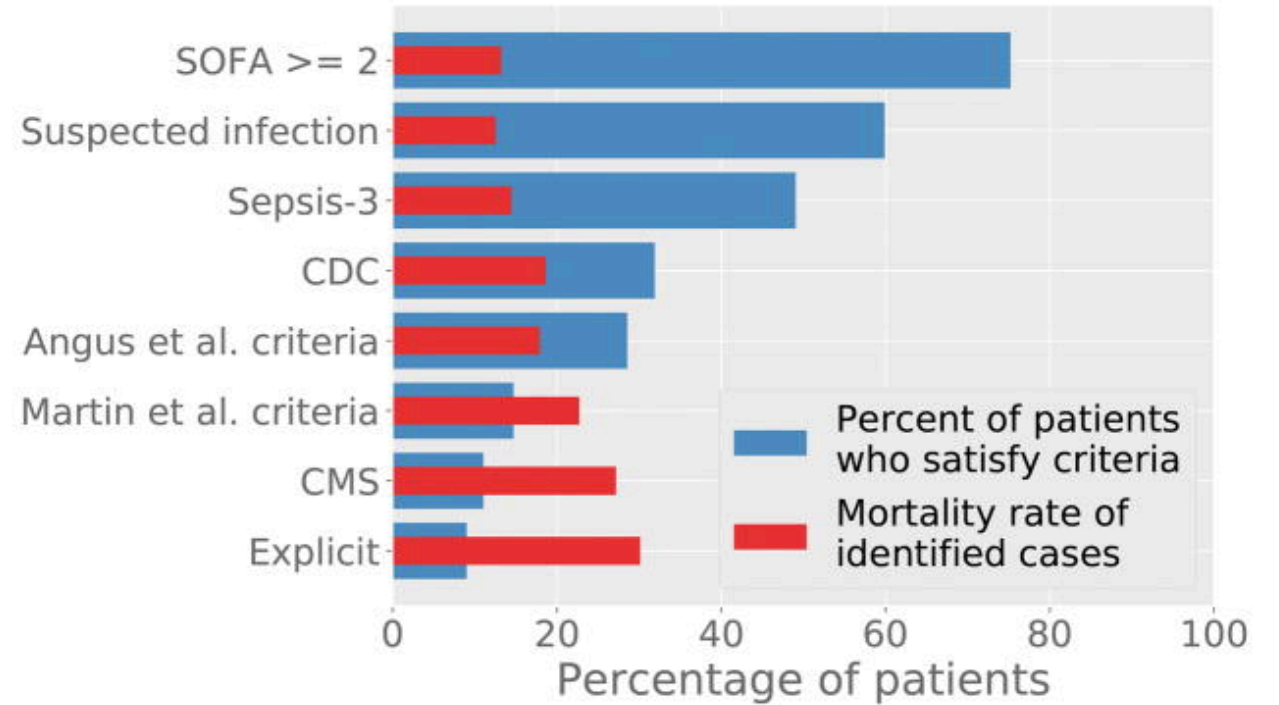
Summary

Background

Previous studies in medical imaging have shown disparate abilities of artificial intelligence (AI) to detect a person's race, yet there is no known correlation for race on medical imaging that would be obvious to human experts when interpreting the images. We aimed

Cohort definitions are important

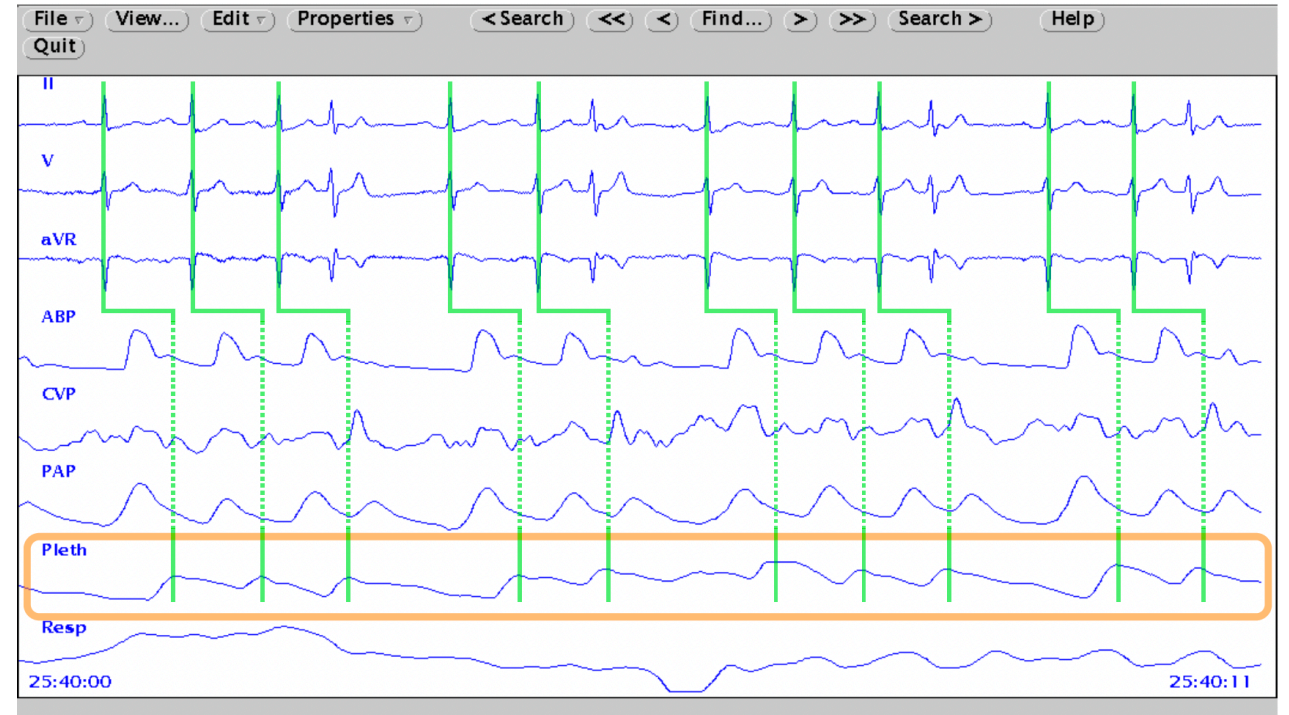
A comparative analysis of sepsis identification methods in an electronic database. Crit Care Med. 2018 Apr;46(4):494-499. doi: 10.1097/CCM.0000000000002965



Different systems use different clocks

Temporal relationships between different streams may be highly informative.

Different clocks on different systems obscure these relationships.



Datasets are often a poor representation of clinical reality

Many public datasets are poor representations of clinical reality.

Promote research that isn't deployable.

Original Investigation | Health Informatics 

Diagnostic Codes in AI Prediction Models and Label Leakage of Same-Admission Clinical Outcomes

Bashar Ramadan, MBBS¹; Ming-Chieh Liu, MS¹; Michael C. Burkhart, PhD¹; [et al](#)

[» Author Affiliations](#) | [Article Information](#)

 [RELATED ARTICLES](#)  [FIGURES](#)  [SUPPLEMENTAL CONTENT](#)

Key Points

Question Are *International Classification of Diseases (ICD)* diagnostic codes, which are only finalized after hospital discharge, associated with inflated performance of artificial intelligence (AI) health care prediction models?

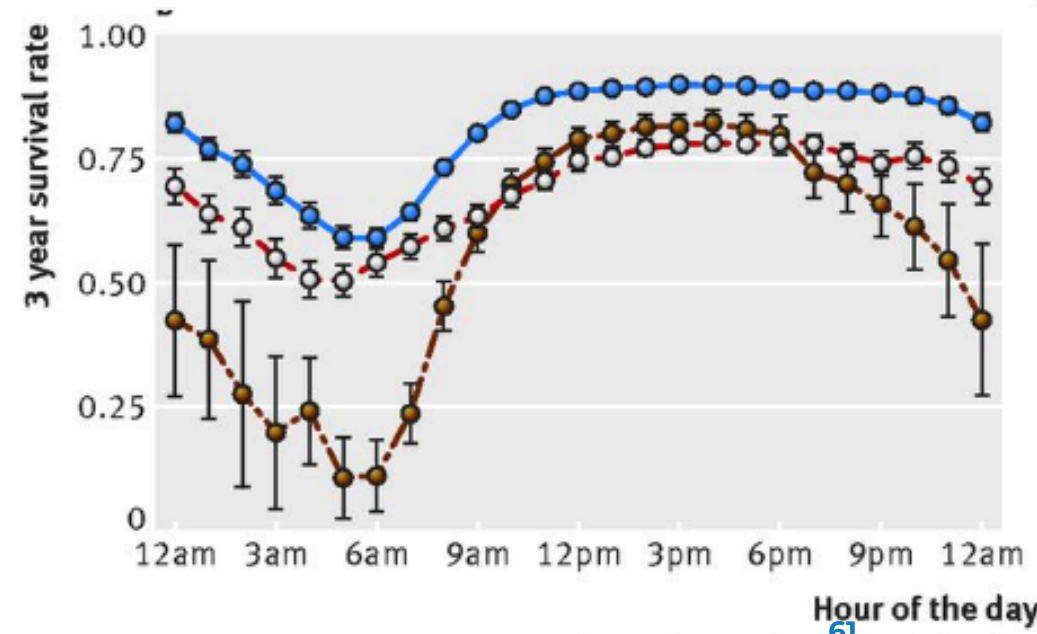
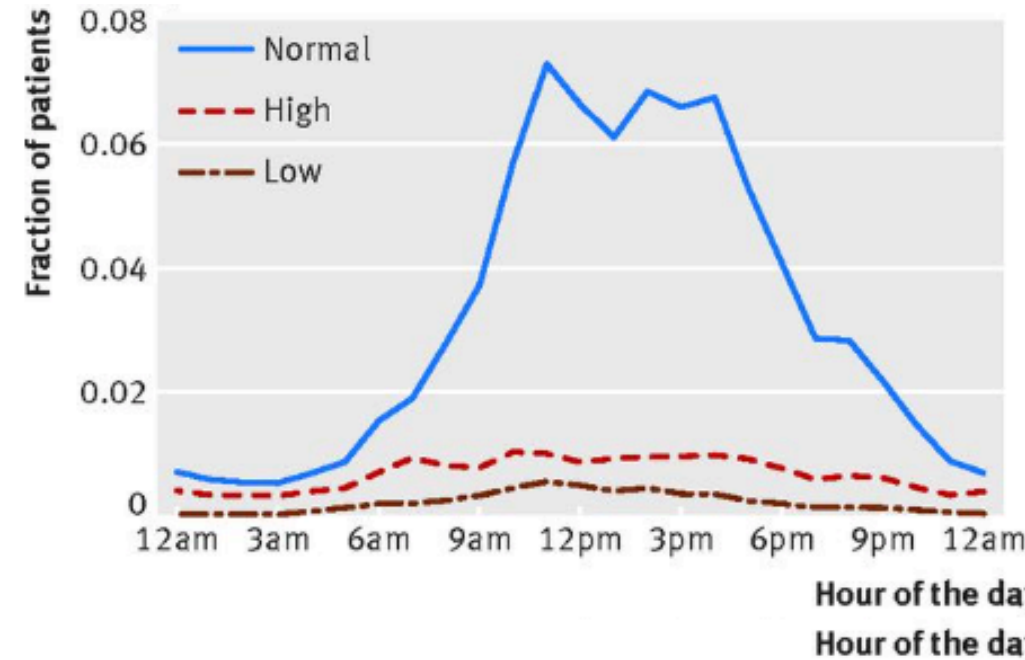
Findings In this prognostic study of 180 640 patients, 40.2% of published AI models trained to predict same-admission outcomes used *ICD* codes as features. Prediction models for inpatient mortality trained on *ICD* codes predicted in-hospital mortality with high accuracy, with the most important codes (eg, brain death, encounter for palliative care) not available in time for clinically useful mortality prediction.

Timing of laboratory tests

"the timing of when laboratory tests were ordered was more accurate than the test results in predicting survival in 118 of 174 tests

patients with normal white blood cell count values taken at 4 am have lower survival (85.4%) than patients with an abnormal measurement at 4 pm

Biases in electronic health record data due to processes within the healthcare system. BMJ 2018;361:k1479



Dataset shift

Occurs when the distribution of data changes between training and deployment

Leads to deterioration of model performance.

<https://www.nejm.org/doi/full/10.1056/NEJMc2104626>

The screenshot shows the top portion of a webpage from The New England Journal of Medicine. At the top, there is a navigation bar with links for 'The New England Journal of Medicine', 'NEJM Evidence', 'NEJM AI', 'NEJM Catalyst', and 'NEJM Journal Watch'. Below this is the journal's logo and name, 'The NEW ENGLAND JOURNAL of MEDICINE', along with dropdown menus for 'CURRENT ISSUE', 'SPECIALTIES', and 'TOPICS'. The article title is 'The Clinician and Dataset Shift in Artificial Intelligence', categorized as 'CORRESPONDENCE'. It includes publication details: 'Published July 14, 2021 | N Engl J Med 2021;385:283-286 | DOI: 10.1056/NEJMc2104626 | VOL. 385 NO. 3' and 'Copyright © 2021'. A row of social media and utility icons (bell, bookmark, copyright, print, quote, share) is visible. The main text begins with 'TO THE EDITOR:' followed by a paragraph: 'Artificial intelligence (AI) systems are now regularly being used in medical settings,¹ although regulatory oversight is inconsistent and undeveloped.^{2,3} Safe deployment of clinical AI requires informed clinician-users, who are generally responsible for identifying and reporting emerging problems. Clinicians may also serve as administrators in governing the use of clinical AI. A natural question follows: are clinicians adequately prepared to identify circumstances in which AI systems fail to perform their intended function reliably?' The next paragraph starts with 'A major driver of AI system malfunction is known as “dataset shift.”^{4,5} Most clinical AI systems today use machine learning, algorithms that leverage statistical methods to learn key'. On the right side of the page, there is a vertical sidebar with icons for information, home, search, and other functions.

Changes in ICD codes

Examined insurance claims for >18m people in the US from 2010 to 2017

Transition from ICD-9 to ICD-10 led to instantaneous +/- of >20% in the prevalence of many diagnostic categories.

<https://jamanetwork.com/journals/jamanetworkopen/full>

Original Investigation | Health Policy



Diagnostic Category Prevalence in 3 Classification Systems Across the Transition to the *International Classification of Diseases, Tenth Revision, Clinical Modification*

Randall P. Ellis, PhD¹; Heather E. Hsu, MD, MPH²; Chenlu Song, MA¹; [et al](#)

» [Author Affiliations](#) | [Article Information](#)

☰ [RELATED ARTICLES](#) [FIGURES](#) [SUPPLEMENTAL CONTENT](#)

Key Points

Question Was the transition from *International Classification of Diseases, Ninth Revision, Clinical Modification (ICD-9-CM)* to the *Tenth Revision (ICD-10-CM)* in October 2015 associated with changes in diagnostic category prevalence when diagnoses are grouped by classification system?

Findings This interrupted time series analysis and cross-sectional study examined insurance claims for more than 18 million privately insured adults and children in the US from 2010 to 2017 and found instantaneous increases or decreases of 20% or more associated with the *ICD-10-CM* transition for nearly 1 in 6 (16%) diagnostic categories

Who is missing entirely

Perhaps the most significant bias is in the patients who **never enter the data**.

- Under-served populations under-appear — or appear only at crisis.
- Missingness is **informative**: absence reflects access, not health.
- A model trained on who showed up encodes who **didn't**.

→ Days 3 & 5: equity

→ data fusion (next)

Frontier — "more data" can make it worse

- Combining sites is **data fusion**, not addition: each source carries its **own selection mechanism**.
- Pooling without modeling those mechanisms can **compound** bias — and *worsen* transport to a new site.

The instinct "just get more data" is right only if you know what each source selected for.

↻ Day 1: more sites ≠ more general

→ Day 5: transport

Data answers the question it was generated for

- Spending and utilization data reflect the **health system**, not just health.
- There is almost always a **mismatch** between what the data *measures* and what you *want to know*.
- Name that gap explicitly — before it silently becomes your result.

AHLI HEALTH AI SUMMER CAMP 2026
BACK TO YOUR OWN PROJECTS

Interrogating your own dataset

Four questions to ask of your data:

- 1 **Where did it come from** — what process generated it, for what purpose?
- 2 **Who is missing** — who never enters, who is under-measured?
- 3 **What does it silently encode** — coding, workflow, system incentives?
- 4 **Does it support your Day 1 problem** — or quietly change it?

The gap check

Does the data that *exists* actually support the problem you posed yesterday?

- If **yes** — say how, specifically.
- If **no** — you must **revise the problem, the data plan, or both.**

A problem statement without a data spec is still a wish.

Your data specification

Interrogate the data your project depends on.

Today's deliverable: Data specification sheet.

Consider **what data exists**, the **generating process**, your **access plan**, and the **anticipated provenance-driven bias**.

Next:

- 11:00am: small group, interrogate your data
- 1:30pm: Sage Bionetworks discuss data beyond the EHR