



UNIVERSITY  
OF  
CALIFORNIA

# Precisely Practicing Medicine from 700 Trillion Points of University of California Health Data

**Atul Butte, MD, PhD**

Chief Data Scientist, University of California Health (UC Health)

Priscilla Chan and Mark Zuckerberg Distinguished Professor

Director, Bakar Computational Health Sciences Institute, UCSF

atul.butte@ucsf.edu ▪ @atulbutte

# Conflicts of Interest

- Scientific founder and advisory board membership
  - Genstruct
  - NuMedii
  - Personalis
  - Carmenta
- Honoraria for talks
  - Lilly
  - Pfizer
  - Siemens
  - Bristol Myers Squibb
  - AstraZeneca
  - Roche
  - Genentech
  - Warburg Pincus
  - CRG
  - AbbVie
  - Westat
- Past or present consultancy
  - Lilly
  - Johnson and Johnson
  - Roche
  - NuMedii
  - Genstruct
  - Tercica
- Ecoeos
  - Helix
  - Ansh Labs
  - uBiome
  - Prevendia
  - Samsung
  - Assay Depot
  - Regeneron
  - Verinata
  - Pathway Diagnostics
  - Geisinger Health
  - Covance
  - Wilson Sonsini Goodrich & Rosati
  - Orrick
  - 10X Genomics
  - GNS Healthcare
  - Gerson Lehman Group
  - Coatue Management
- Corporate Relationships
  - Northrop Grumman
  - Genentech
  - Optum
  - Aptalis
  - Allergan
  - Astellas
  - Thomson Reuters
- Intel
  - SAP
  - SV Angel
  - Progenity
  - Illumina
- Speakers' bureau
  - None
- Companies started by students
  - Carmenta
  - Serendipity
  - Stimulomics
  - NunaHealth
  - Praedicat
  - MyTime
  - Flipora
  - Tumbl.in
  - Polyglot
  - Iota Health
  - Ongevity Health

# University of California

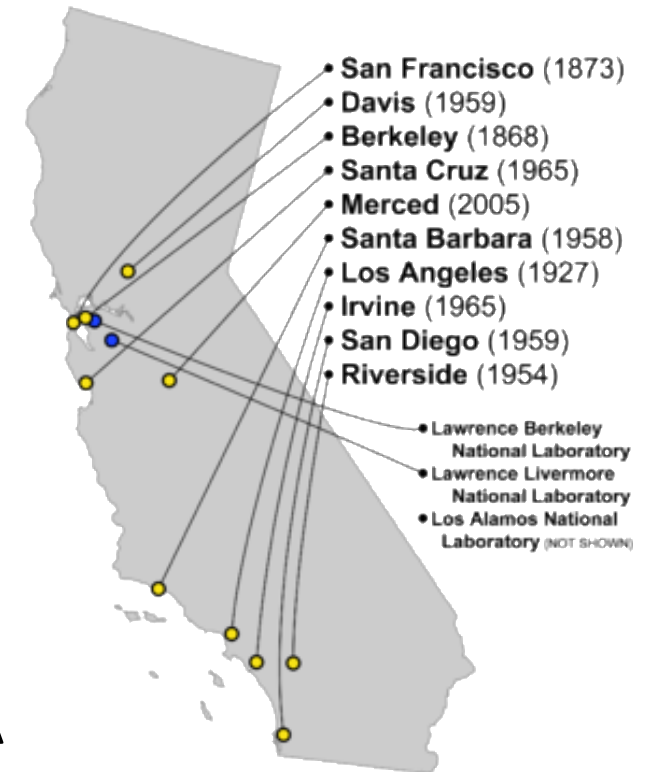
- 10 campuses and 3 national labs
- ~200,000 employees, ~250,000 students/yr



UC Health

## UC Health

- 18 health professional schools (6 med schools)
- Train half the medical students and residents in California
- ~\$2 billion NIH funding
- \$13+ billion clinical operating revenue
- 5000 faculty physicians, 12000 nurses
- UCSF and UCLA are in US News top 10
- 5 NCI Comprehensive Cancer Centers, 5 NIH CTSA
- IRB reliance, centralized contracting



# UC Health, United Healthcare Form New ACO & Clinically Integrated Networks

by Staff Writer

📅 10/03/2016



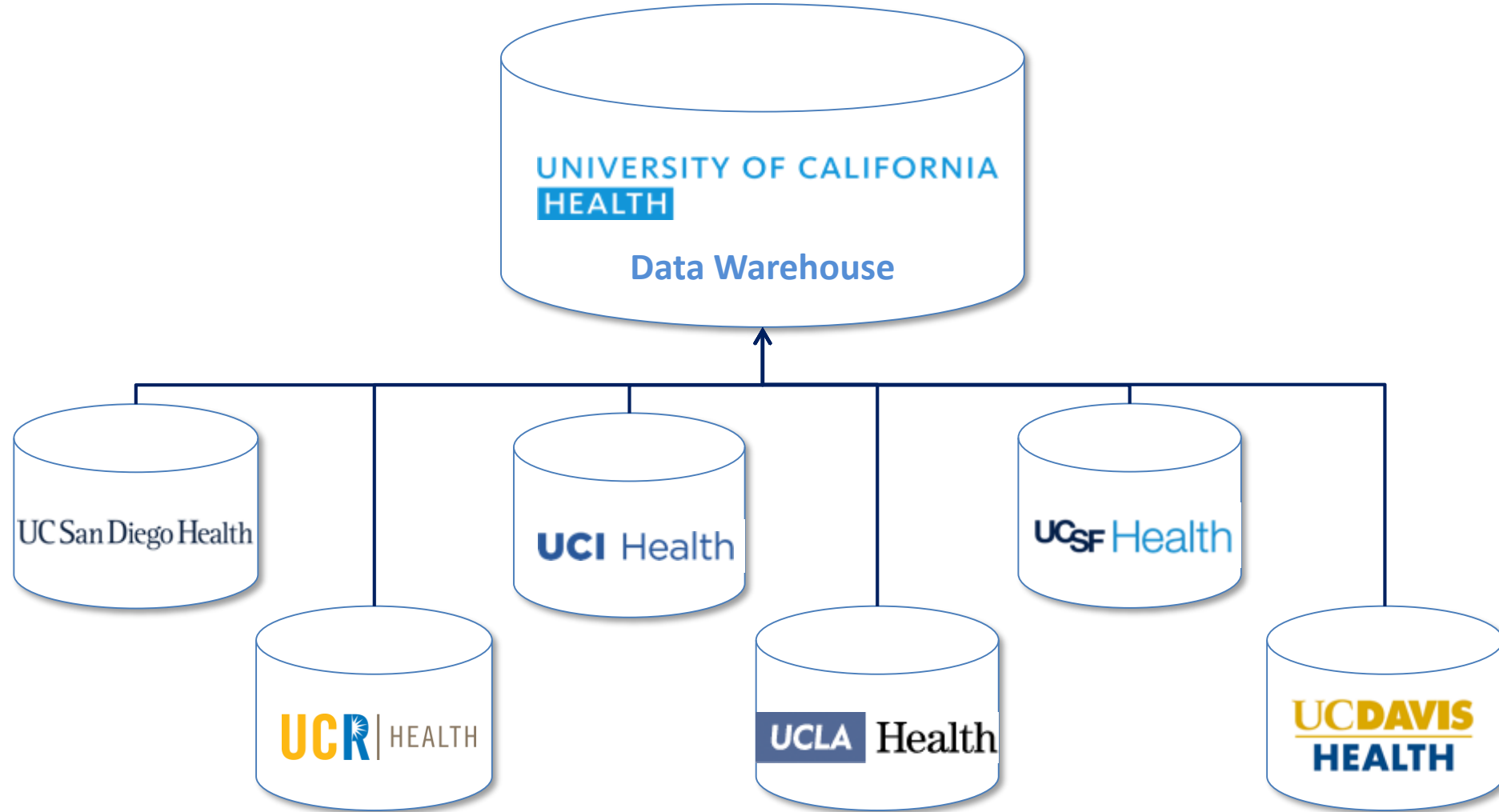
0 Comments



## UC Health

The University of California and UnitedHealth Group are teaming up to form a new accountable care organization (ACO) and clinically integrated network. As part of the 10-year strategic relationship, UC Health's five academic medical centers will expand use of Optum's clinically integrated network services and advanced data analytics services.

# *Combining healthcare data from across the six University of California medical schools and systems*



# The University of California has an incredible view of the medical system

- Combined EHR data from UCSF, UCLA, UC Irvine, UC Davis, UC San Diego, and UC Riverside
  - 15 million patients treated over the past 15 years
- Central database built using OMOP (not Epic) as a data backend
  - Structured data from 2012 to the present day
  - **5.7 million patients with “modern” data**
  - 228M encounters, 96M procedures, 593M med orders, 640M diagnosis codes, 2.3B lab tests and vital signs
  - “Everything from Tylenol to CAR-T cells...”
  - OSHPD data, pathology and radiology text elements, death index
  - Claims data from our self-funded plans now included
  - Continually harmonizing elements
- Quality and performance dashboards

# Many operational teams within UC Health now using and benefitting from the UC Health Data Warehouse, saving \$millions



Central tools to  
improve quality  
of care



Managing costs  
in our self-  
funded health  
plans



Decreasing  
expenses in our  
self-funded  
health plans



Central  
management of  
primary care  
patients



# University of California Cancer Consortium Takes on California's \$14 Billion Killer

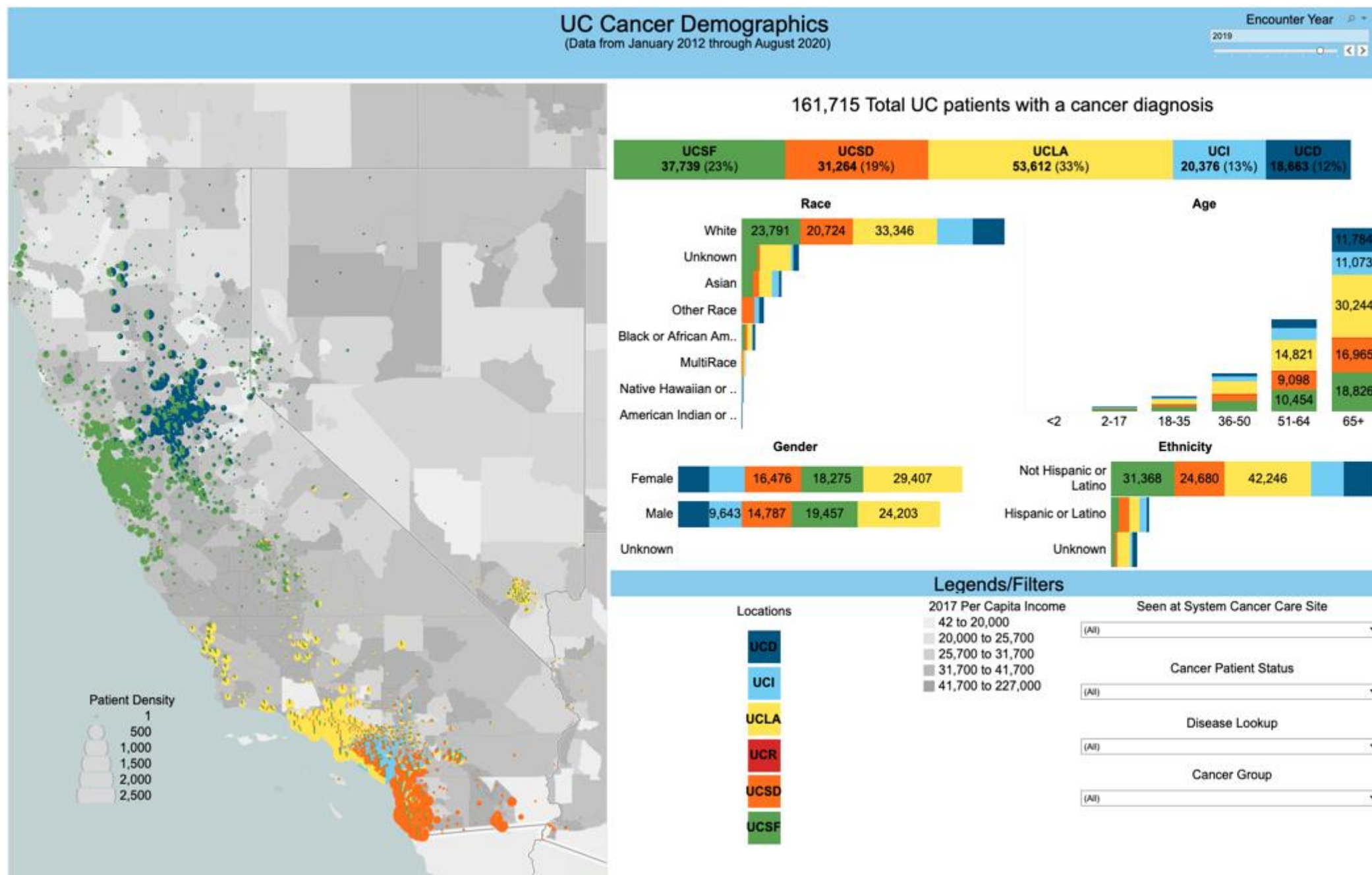
By [Elizabeth Fernandez](#) on September 11, 2017



University of California President Janet Napolitano announced the formation of a UC Cancer Consortium to include five of the nation's leading academic cancer centers during a press conference in Genentech Hall at Mission Bay. *Photo by Susan Merrell*



# 161,715 cancer patients seen in 2019 across UC Health

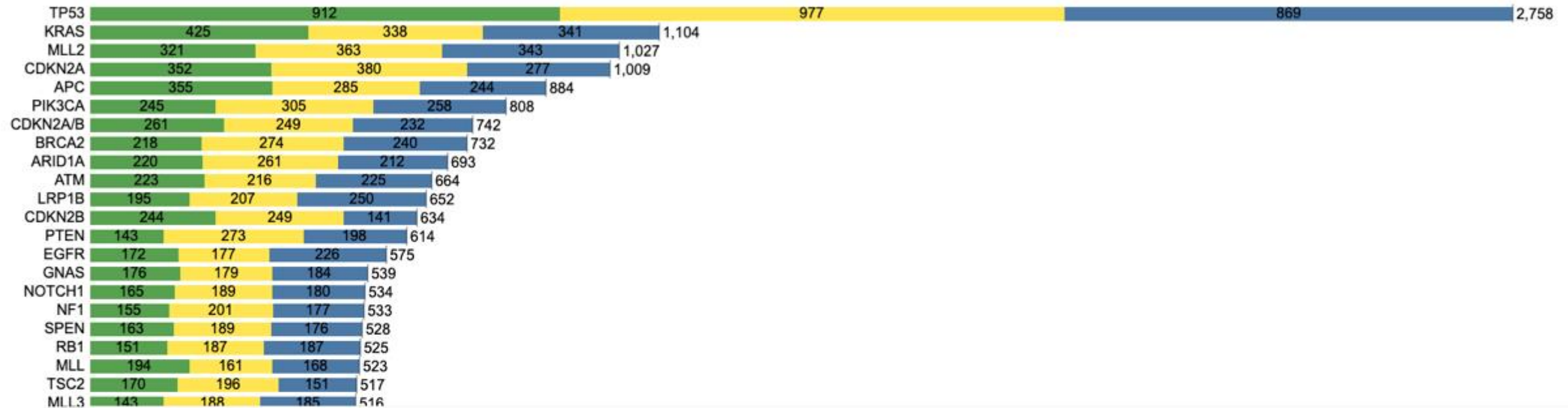


# Foundation Medicine cancer genomic reports centralized and parsed

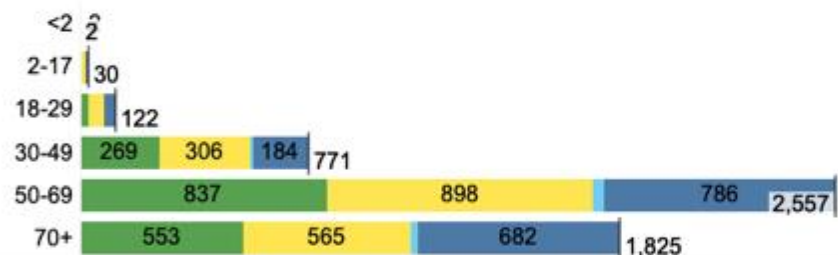
Disease Selection

(All)

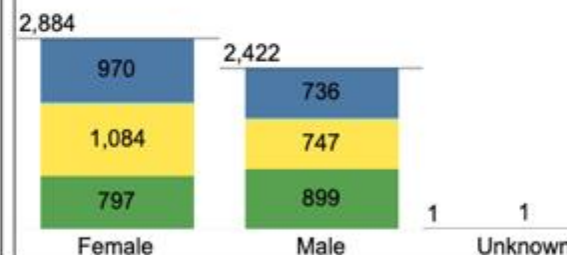
Gene occurrence



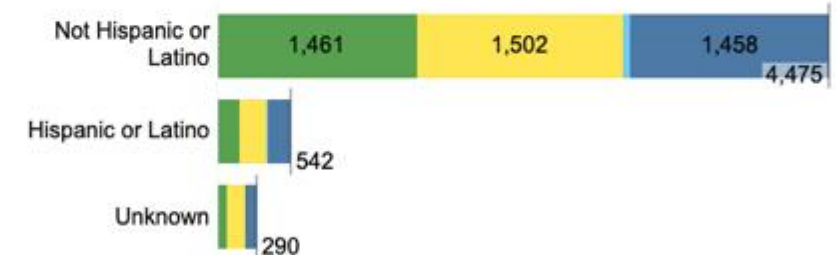
Age



Gender



Ethnicity



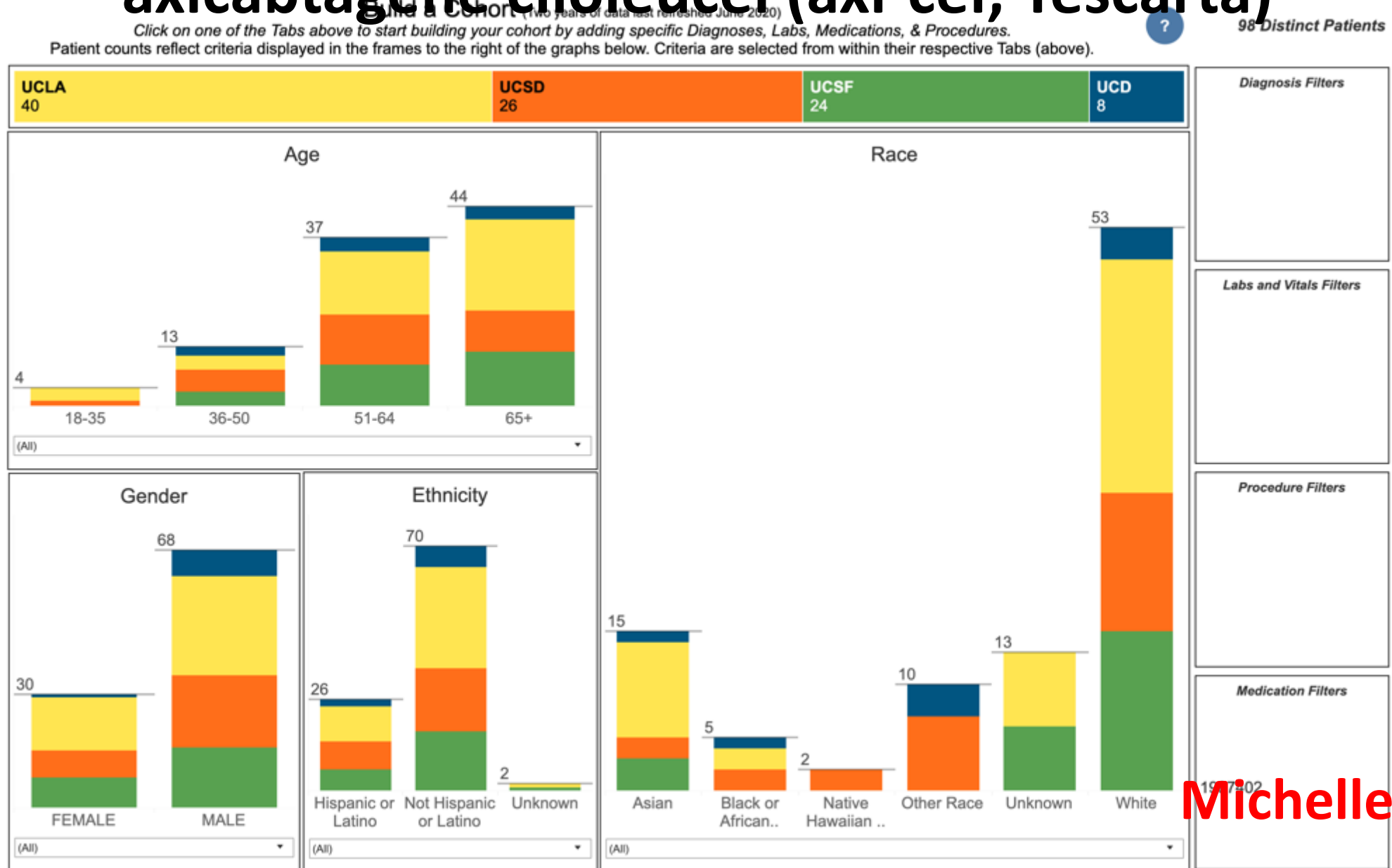
Race

Smoking Status

# Top 10 Drug Charges across UC Health



# Nearly 100 patients treated so far with axicabtagene ciloleucel (axi-cel; Yescarta)



Michelle Wang

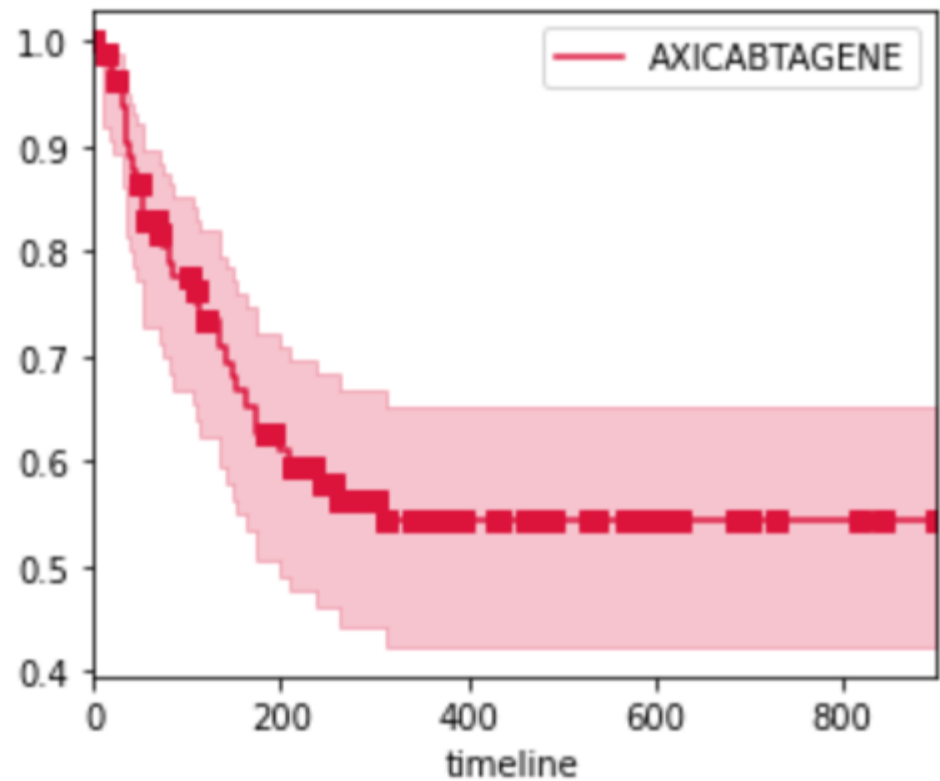






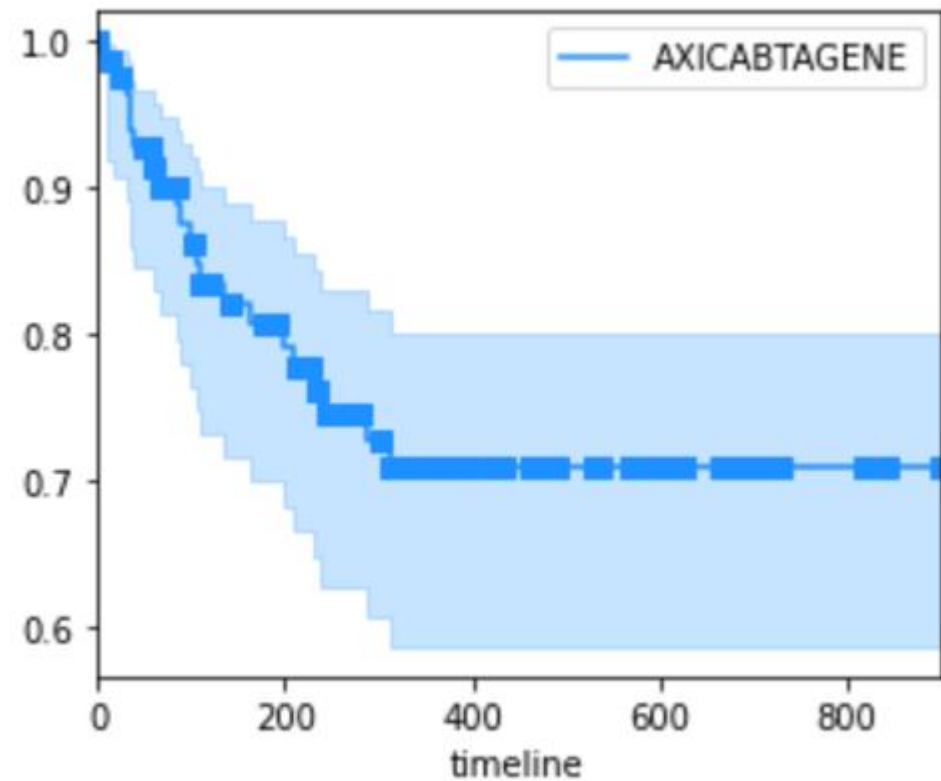
# PFS and OS at 12 month

Progression Free Survival



At risk  
KM\_estimate 85 59 43 32 22 13 9 5 3

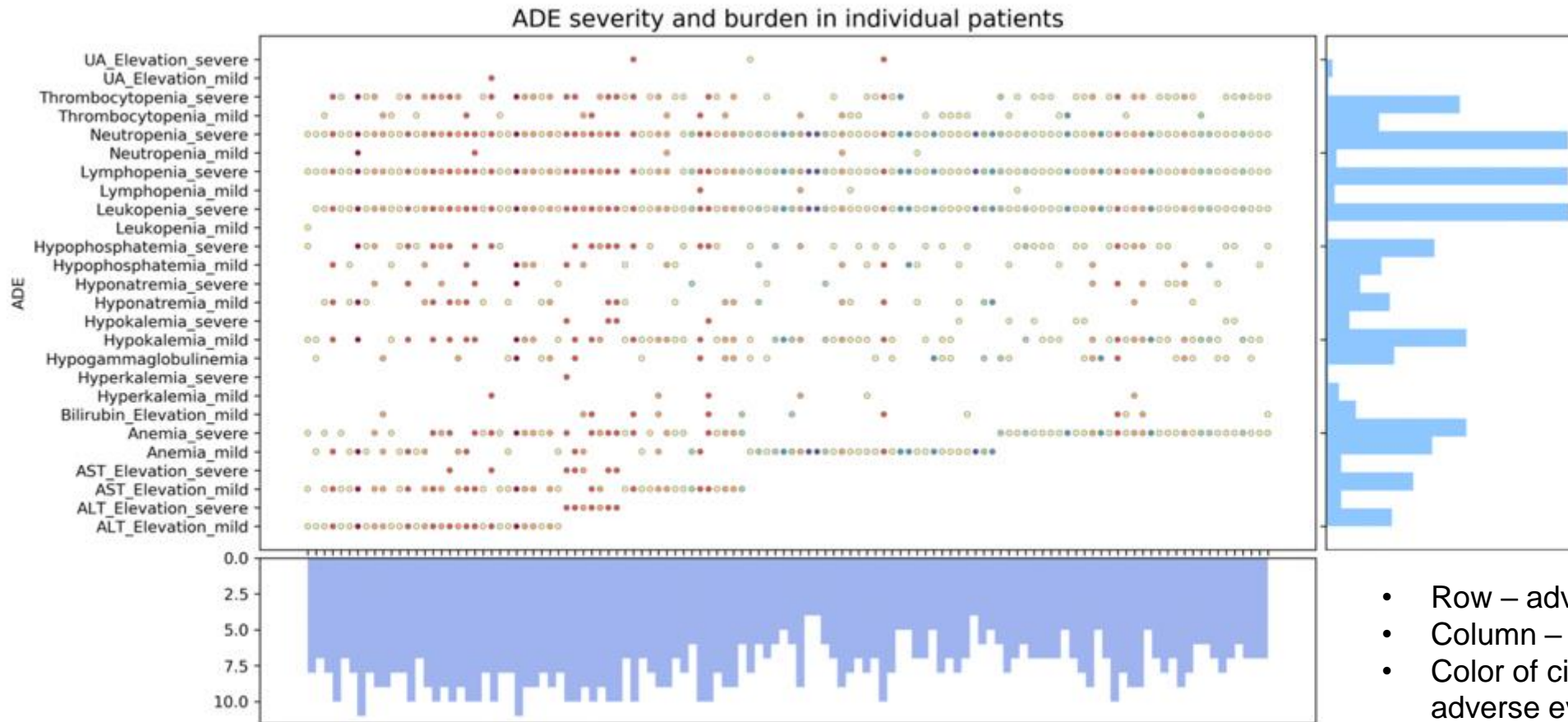
Overall Survival



At risk  
KM\_estimate 85 65 54 41 26 16 11 6 3

Michelle Wang

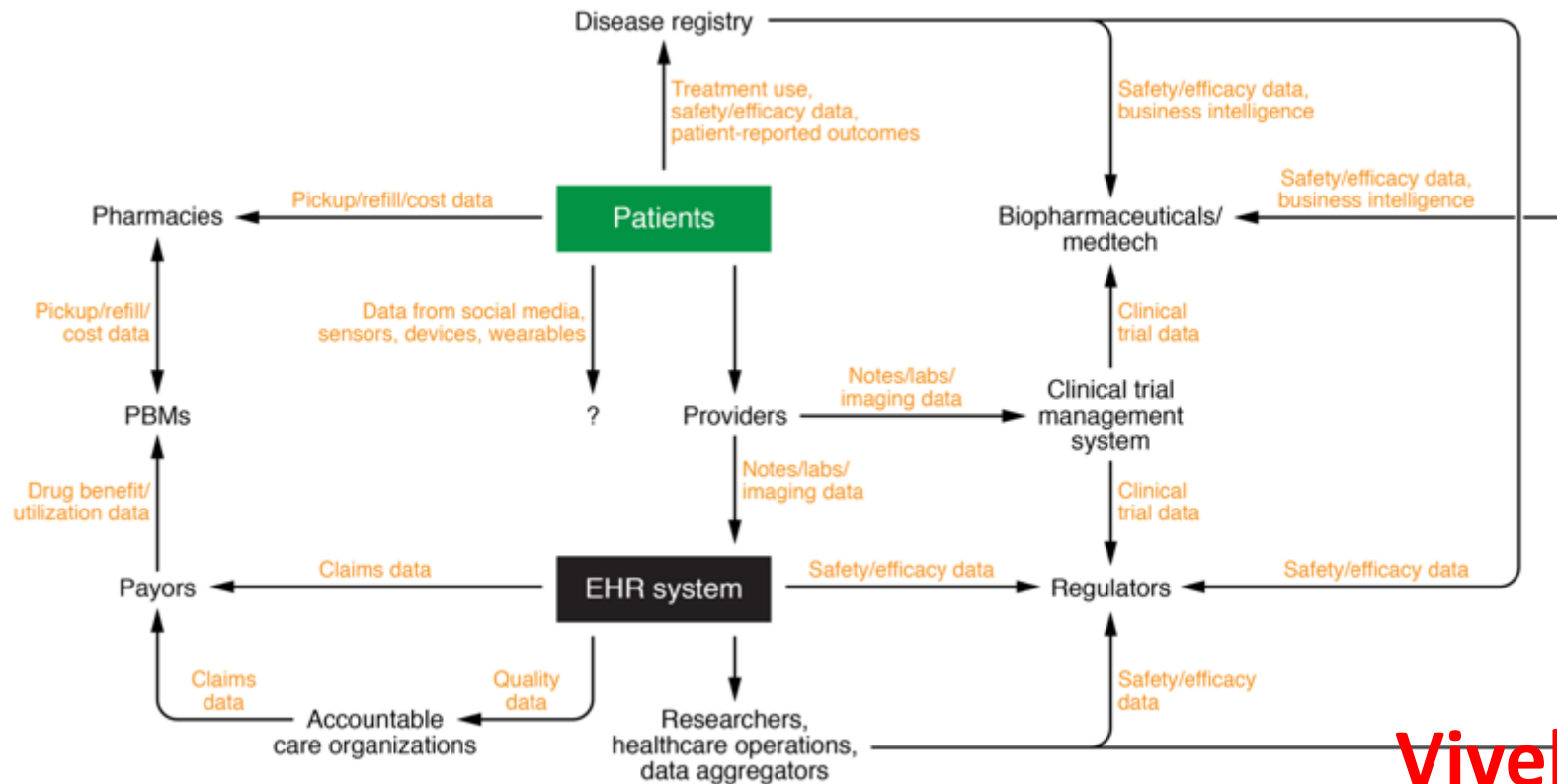
# Overview of ADE burden in individual patients



**Michelle Wang**

# Announcing a new publication on Real World Data

- Journal of Clinical Investigation (JCI) February 2020
  - Opportunities and Challenges in Using Real-World Data for Healthcare
  - Special issue on Big Data's Future in Medicine [bit.ly/JCIbigdata](https://bit.ly/JCIbigdata)



**Vivek Rudrapatna**  
**[bit.ly/JCIbigdata](https://bit.ly/JCIbigdata)**

# Twenty-one uses for Real World Data

## **Post-approval safety**

- Updating side effect rates
- Discovering novel side effects

## **Supporting regulatory approval**

- Single-arm experimental trials
- "Digital approvals"
- Biosimilar development

## **Informing clinical trials design**

- Better patient selection
- Trimming the trials: more efficient data collection

## **Continually establishing efficacy**

- Assessing the efficacy-effectiveness gap
- Searching for efficacy in specific populations
- Effect modifiers and precision medicine
- Long-term, post trial outcomes

## **Comparative effectiveness**

- Integrating costs with comparative effectiveness
- Understanding effects of pharmacy practices on healthcare utilization
- Studying novel on-label pharmaceuticals versus older off-label drugs

## **Studying the practice of medicine**

- Quality of practice, medical errors
- Standardizing care and care delivery
- The effect of payors on medical care
- Are new-generation diagnostics improving outcomes?

## **Data driven decision support**

- Clinical-Decision Support: the Provider Perspective
- Clinical-Decision Support: the Patient Perspective
- Clinical-Decision Support: the Community Perspective

**Vivek Rudrapatna**  
**[bit.ly/JCIbigdata](https://bit.ly/JCIbigdata)**

## **UCLA Health**

Mohammed Mahbouba  
Albert Duntugan  
Jay Shah  
Vajra Kasturi  
Danielle Belmontez  
Michael Swinford (Azure)  
Scott Bailey (Azure)  
Andrew Weaver (Azure)

## **UCI Health**

Lisa Dahm  
Ayan Patel  
Charles Wilson  
Aiden Barin  
Kathy Pickell  
David Gonzalez  
Lattice Armstead  
Tim Hayes

## **UCDAVIS HEALTH**

Kent Anderson  
Doug Berman  
Steve Covington  
Jeff Evoy  
Raj Sankala  
Hemanth Tatiparthi  
Brian Paciotti

## **UCSF Health**

David Dobbs  
Rick Larsen  
Nelson Lee  
Brian Chan  
Oksana Gologorskaya  
Rao Venigalla

# **The Team**

## **UC San Diego Health**

Josh Glandorf  
Jennifer Holland  
Eugene Lee  
Travis Mitchelar  
Peter Ryan

## **The CIO Team**

Mike Pfeffer\* – UCLA  
Chris Longhurst – UCSD  
Joe Bengfort – UCSF  
Chuck Podesta – UCI  
John Cook - UCD

\*UCHDW Sponsoring CIO

## **UC HEALTH**

Atul Butte  
Tom Andriola  
Liz Engel  
Jack Stobo  
Leslie Yuan





# Enabling UC researchers and patients to go beyond... machine learning in a safe, respectful, fair, equitable way in medicine

Alvin Rajkumar<sup>1,2</sup>, Eyal Oren<sup>1</sup>, Kai Chen<sup>1</sup>, Andrew M. Dai<sup>1</sup>, Nissan Hajaj<sup>1</sup>, Michaela Hardt<sup>1</sup>, Peter J. Liu<sup>1</sup>, Xiaobing Liu<sup>1</sup>, Jake Marcus<sup>1</sup>, Mimi Sun<sup>1</sup>, Patrik Sundberg<sup>1</sup>, Hector Yee<sup>1</sup>, Kun Zhang<sup>1</sup>, Yi Zhang<sup>1</sup>, Gerardo Flores<sup>1</sup>, Gavin E. Duggan<sup>1</sup>, Jamie Irvine<sup>1</sup>, Quoc Le<sup>1</sup>, Kurt Litsch<sup>1</sup>, Alexander Mossin<sup>1</sup>, Justin Tansuwan<sup>1</sup>, De Wang<sup>1</sup>, James Wexler<sup>1</sup>, Jimbo Wilson<sup>1</sup>, Dana Ludwig<sup>2</sup>, Samuel L. Volchenbourn<sup>1</sup>, Katherine Chou<sup>1</sup>, Michael Pearson<sup>1</sup>, Srinivasan Madabushi<sup>1</sup>, Nigam H. Shah<sup>4</sup>, Atul J. Butte<sup>2</sup>, Michael D. Howell<sup>1</sup>, Claire Cui<sup>1</sup>, Greg S. Corrado<sup>1</sup> and Jeffrey Dean<sup>1</sup>

Predictive modeling with electronic health record (EHR) data is anticipated to drive personalized medicine and improve healthcare quality. Constructing predictive statistical models typically requires extraction of curated predictor variables from normalized EHR data, a labor-intensive process that discards the vast majority of information in each patient's record. We propose a representation of patients' entire raw EHR records based on the Fast Healthcare Interoperability Resources (FHIR) format. We demonstrate that deep learning methods using this representation are capable of accurately predicting multiple medical events from multiple centers without site-specific data harmonization. We validated our approach using de-identified EHR data from two US academic medical centers with 216,221 adult patients hospitalized for at least 24 h. In the sequential format we propose, this volume of EHR data unrolled into a total of 46,864,534,945 data points, including clinical notes. Deep learning models achieved high accuracy for tasks such as predicting: in-hospital mortality (area under the receiver operator curve [AUROC] across sites 0.93–0.94), 30-day unplanned readmission (AUROC 0.75–0.76), prolonged length of stay (AUROC 0.85–0.86), and all of a patient's final discharge diagnoses (frequency-weighted AUROC 0.90). These models outperformed traditional, clinically-used predictive models in all cases. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. In a case study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

npj Digital Medicine (2018)1:18 | doi:10.1038/s41746-018-0029-1

## INTRODUCTION

The promise of digital medicine stems in part from the hope that, by digitizing health data, we might more easily leverage computer information systems to understand and improve care. In fact, routinely collected patient healthcare data are now approaching the genomic scale in volume and complexity.<sup>1</sup> Unfortunately, most of this information is not yet used in the sorts of predictive statistical models clinicians might use to improve care delivery. It is widely suspected that use of such efforts, if successful, could provide major benefits not only for patient safety and quality but also in reducing healthcare costs.<sup>2–6</sup>

In spite of the richness and potential of available data, scaling the development of predictive models is difficult because, for traditional predictive modeling techniques, each outcome to be predicted requires the creation of a custom dataset with specific variables.<sup>7</sup> It is widely held that 80% of the effort in an analytic model is preprocessing, merging, customizing, and cleaning datasets,<sup>8,9</sup> not analyzing them for insights. This profoundly limits the scalability of predictive models.

Another challenge is that the number of potential predictor variables in the electronic health record (EHR) may easily number in the thousands, particularly if free-text notes from doctors,

nurses, and other providers are included. Traditional modeling approaches have dealt with this complexity simply by choosing very limited number of commonly collected variables to consider. This is problematic because the resulting models may produce imprecise predictions: false-positive predictions can overwhelm physicians, nurses, and other providers with false alarms and concomitant alert fatigue,<sup>10</sup> which the Joint Commission identifies as a national patient safety priority in 2014.<sup>11</sup> False-negative predictions can miss significant numbers of clinically important events, leading to poor clinical outcomes.<sup>11,12</sup> Incorporating the entire EHR, including clinicians' free-text notes, offers some hope of overcoming these shortcomings but is unwieldy for most predictive modeling techniques.

Recent developments in deep learning and artificial neural networks may allow us to address many of these challenges and unlock the information in the EHR. Deep learning emerged as the preferred machine learning approach in machine perceptual problems ranging from computer vision to speech recognition but has more recently proven useful in natural language processing, sequence prediction, and mixed modality data settings.<sup>13–17</sup> These systems are known for their ability to handle large volumes of relatively messy data, including errors in labels

## Patient Timeline



<sup>1</sup>Google Inc, Mountain View, CA, USA; <sup>2</sup>University of California, San Francisco, San Francisco, CA, USA; <sup>3</sup>University of Chicago Medicine, Chicago, IL, USA and <sup>4</sup>Stanford University, Stanford, CA, USA  
Correspondence: Alvin Rajkumar (alvinrajkumar@google.com)  
These authors contributed equally: Alvin Rajkumar, Eyal Oren

Received: 26 January 2018 Revised: 14 March 2018 Accepted: 26 March 2018  
Published online: 08 May 2018

# Predicting the future state of a patient with Rheumatoid Arthritis

JAMA  
Network | Open



Original Investigation | Health Informatics

## Assessment of a Deep Learning Model Based on Electronic Health Record Data to Forecast Clinical Outcomes in Patients With Rheumatoid Arthritis

Beau Norgeot, MS; Benjamin S. Glicksberg, PhD; Laura Trupin, MPH; Dmytro Litulev, PhD; Milena Gianfrancesco, PhD, MPH; Boris Oskotsky, PhD; Gabriela Schmajuk, MD, MSc; Jinoos Yazdany, MD, MPH; Atul J. Butte, MD, PhD

### Abstract

**IMPORTANCE** Knowing the future condition of a patient would enable a physician to customize current therapeutic options to prevent disease worsening, but predicting that future condition requires sophisticated modeling and information. If artificial intelligence models were capable of forecasting future patient outcomes, they could be used to aid practitioners and patients in prognosticating outcomes or simulating potential outcomes under different treatment scenarios.

**OBJECTIVE** To assess the ability of an artificial intelligence system to prognosticate the state of disease activity of patients with rheumatoid arthritis (RA) at their next clinical visit.

**DESIGN, SETTING, AND PARTICIPANTS** This prognostic study included 820 patients with RA from rheumatology clinics at 2 distinct health care systems with different electronic health record platforms: a university hospital (UH) and a public safety-net hospital (SNH). The UH and SNH had substantially different patient populations and treatment patterns. The UH has records on approximately 1 million total patients starting in January 2012. The UH data for this study were accessed on July 1, 2017. The SNH has records on 65 000 unique individuals starting in January 2013. The SNH data for the study were collected on February 27, 2018.

**EXPOSURES** Structured data were extracted from the electronic health record, including exposures (medications), patient demographics, laboratories, and prior measures of disease activity. A longitudinal deep learning model was used to predict disease activity for patients with RA at their next rheumatology clinic visit and to evaluate interhospital performance and model interoperability strategies.

**MAIN OUTCOMES AND MEASURES** Model performance was quantified using the area under the

### Key Points

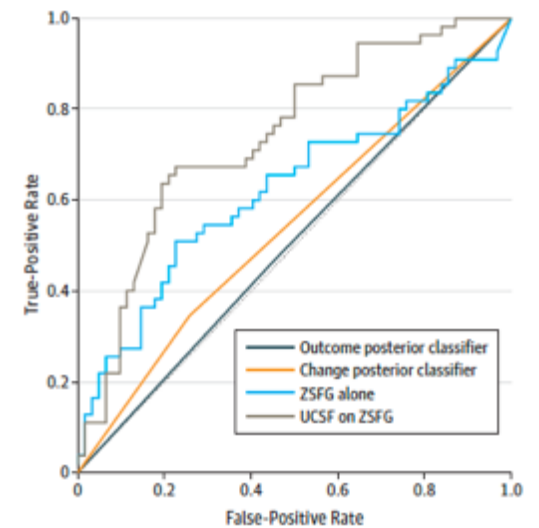
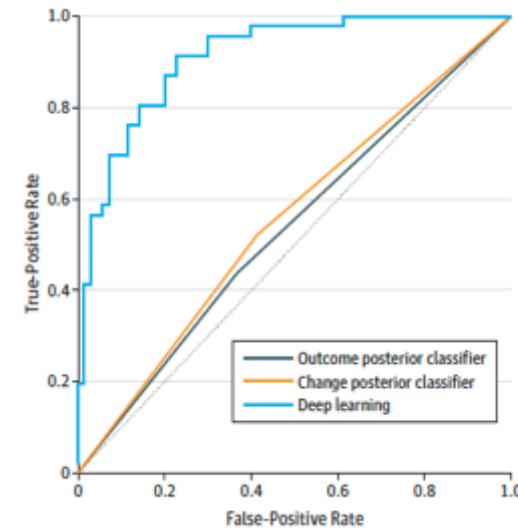
**Question** How accurately can artificial intelligence models prognosticate future patient outcomes for a complex disease, such as rheumatoid arthritis?

**Findings** In this prognostic study of 820 patients with rheumatoid arthritis, a longitudinal deep learning model had strong performance in a test cohort of 116 patients, whereas baselines that used each patient's most recent disease activity score had statistically random performance.

**Meaning** The findings suggest that building accurate models to forecast complex disease outcomes using electronic health records is possible.

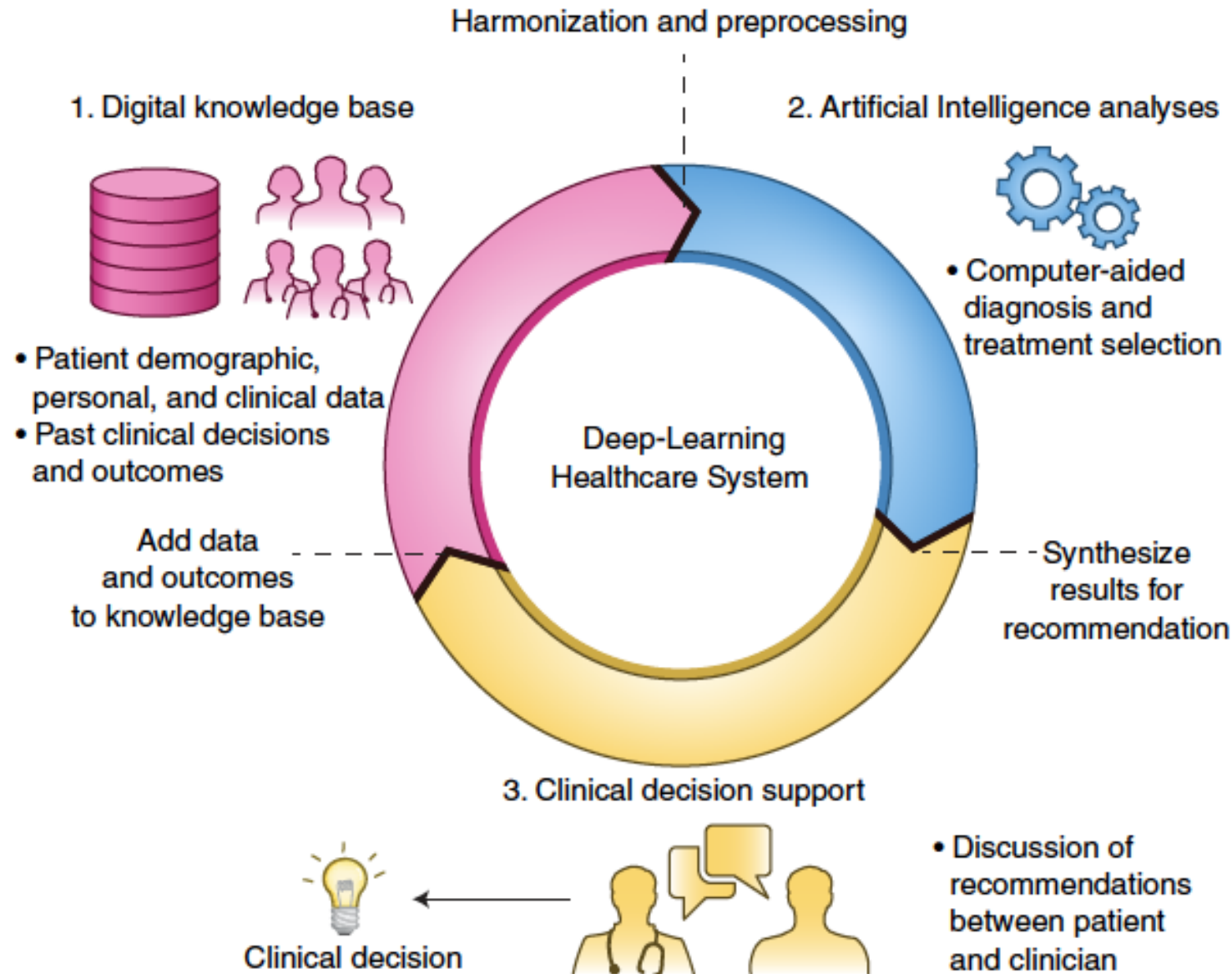
+ Supplemental content

Author affiliations and article information are listed at the end of this article.



Beau Norgeot  
[bit.ly/jamaRA](https://bit.ly/jamaRA)

# A New Deep Learning Healthcare System



**Beau Norgeot**  
**[bit.ly/DLHCare](https://bit.ly/DLHCare)**

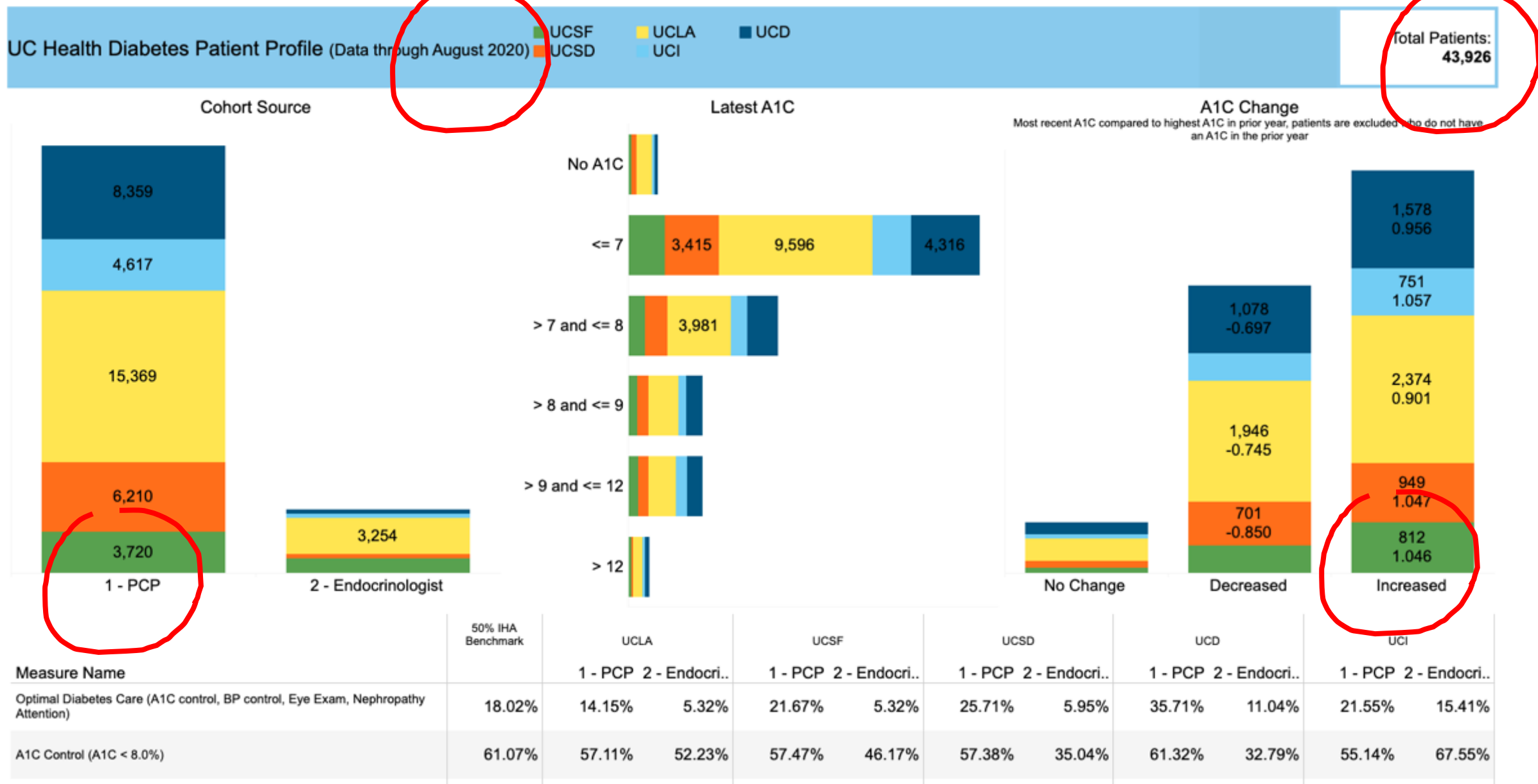


# What could we do with clinical data?

- Clinical researcher at UCLA could run a genome wide association study across UC Health
- Mobile health researcher at UCSD can enable patients to contribute data for research
- Community activist and researcher UC Merced can study environmental factors contributing to health and disease
- Transplant patient at UC Irvine can download all their data across UC Health
- Data scientist at UC Santa Barbara can model development of Alzheimer's disease and build a multi-modal predictor
- App designer at UC Riverside can show patients their choices with chronic disease
- CMO at UCSF can build predictive models for readmission, test, share across UC Health
- AI researcher at UC Berkeley can build deep-learning models for image-based diagnostics
- Health services researcher at UC Davis can build predictive models for drug efficacy, and maybe enable pay-for-performance
- Cancer genomics researcher at UCSC can study all our clinical cancer genomes



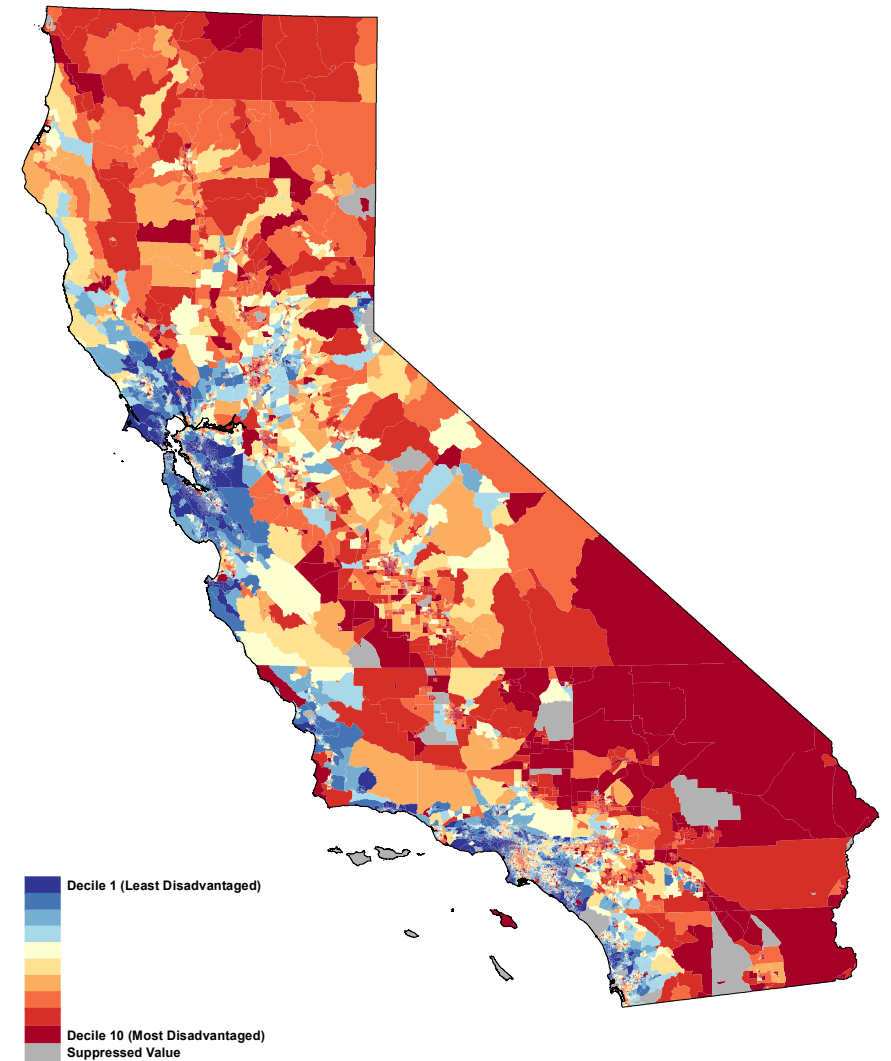
# One dashboard for primary care and specialists covers all 44 thousand UC Health patients with type 2 diabetes



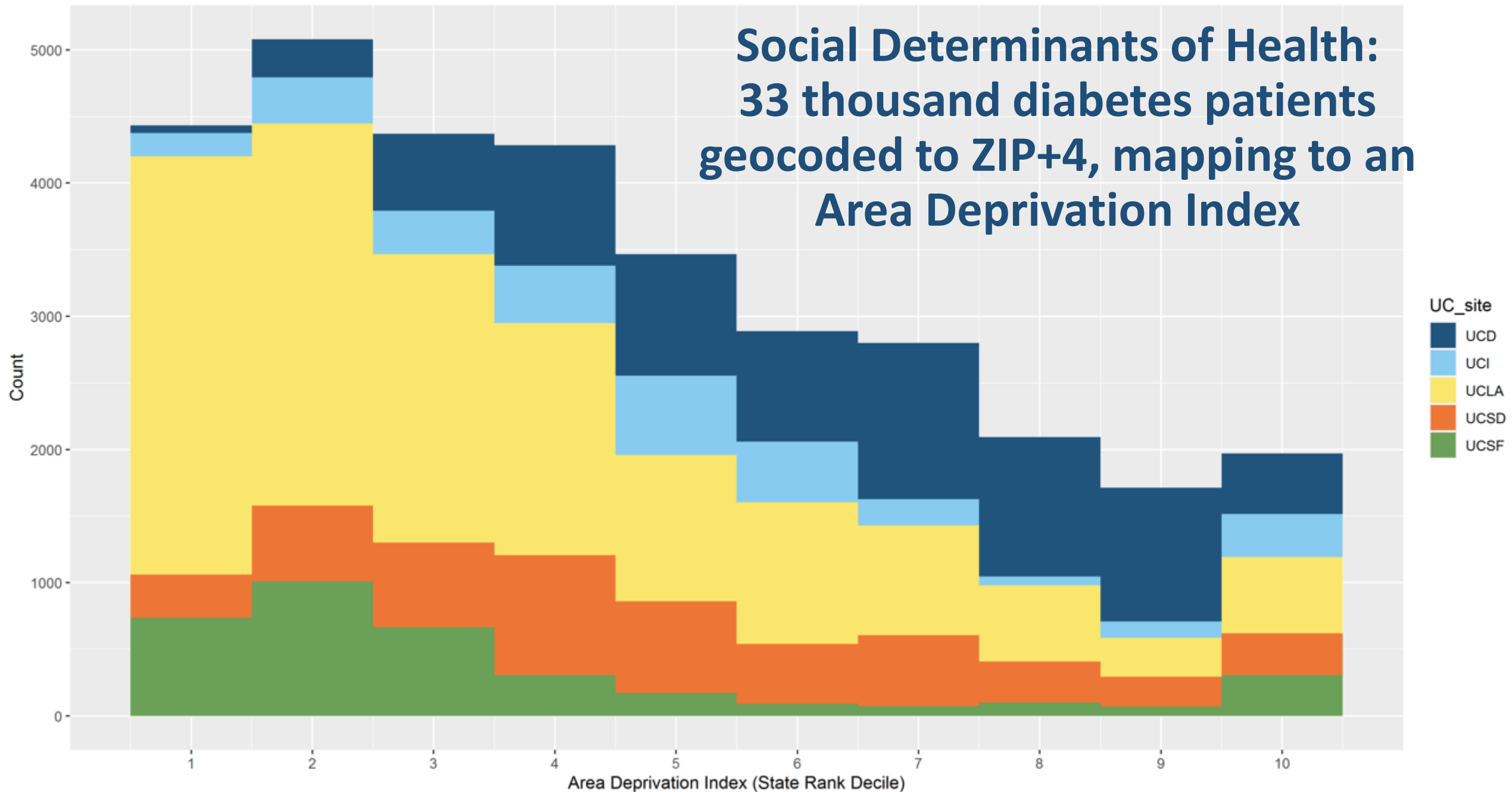
# Social Determinants of Health: Area Deprivation Index (ADI)

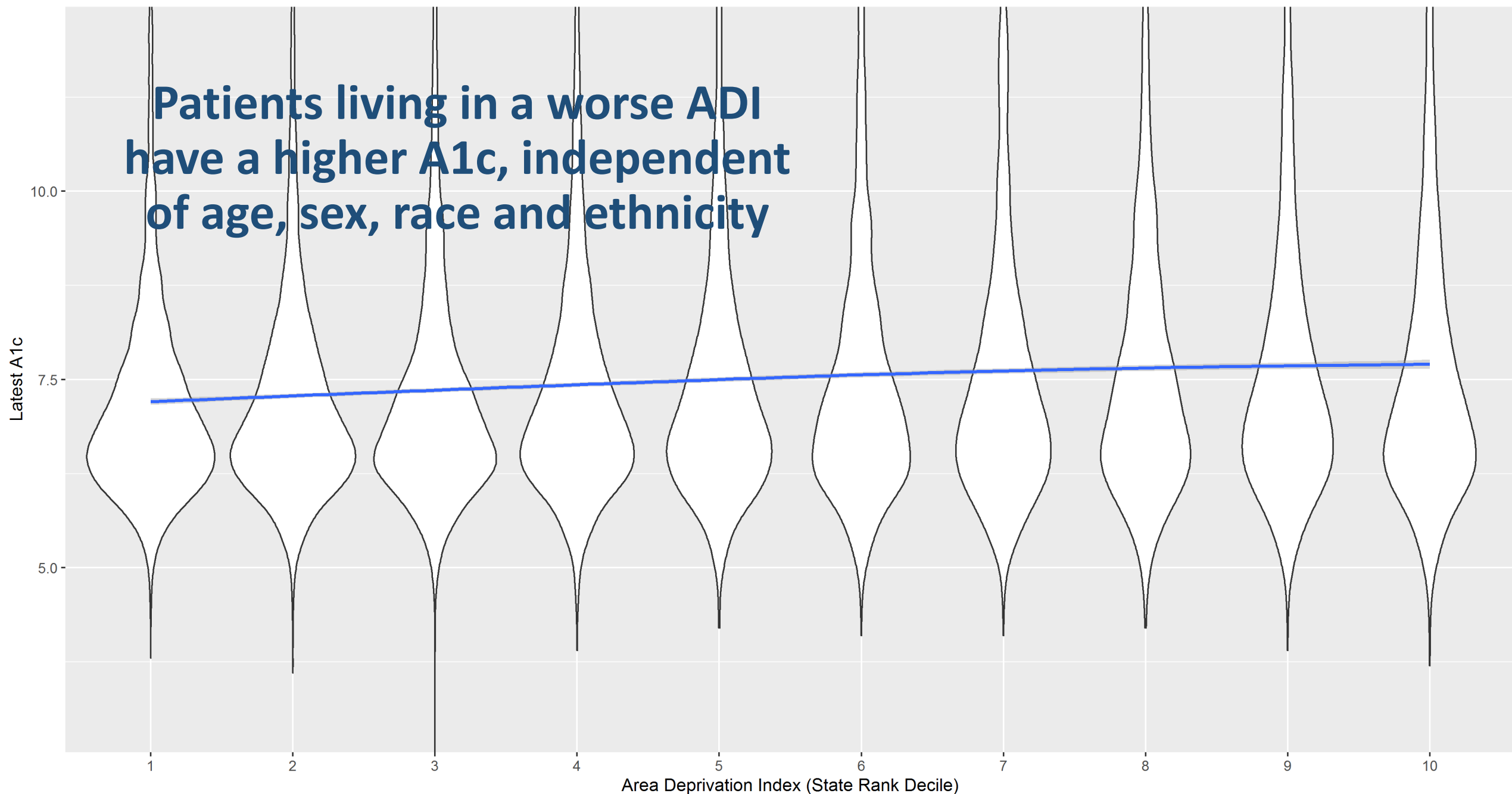
California - 2015 ADI State Rankings

- Estimate socioeconomic status based on income, education, employment, housing quality at the neighborhood level
- Use with race, ethnicity, gender, age to identify health disparities
- Central geocoding capability to link 9-digit zip code for addresses
- We have now geocoded our UC-wide primary care population
- ADI is significantly associated with adverse health outcomes in our patients, in addition to race, ethnicity, and age



# Social Determinants of Health: 33 thousand diabetes patients geocoded to ZIP+4, mapping to an Area Deprivation Index





**Working closely with our pharmacist teams, inappropriate use of IV acetaminophen has significantly dropped over the past 2.5 years**

