

Precisely Practicing Cancer 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
 - Personalis
 - NuMedii
 - Carmenta (Progenity)
 - Genstruct
- Honoraria for talks
 - Lilly
 - Pfizer
 - Siemens
 - Bristol Myers Squibb
 - AstraZeneca
 - Roche
 - Genentech
 - Warburg Pincus
 - CRG
 - AbbVie
 - Westat
- Past or present consultancy
 - Personalis
 - NuMedii
 - Lilly
 - Johnson and Johnson
 - Roche

- Genstruct
- Tercica
- Ecoeos
- Helix
- Ansh Labs
- uBiome
- Prevendia
- Samsung
- Assay Depot
- Regeneron
- Verinata (Illumina)
- Pathway Diagnostics
- Geisinger Health
- Covance
- Wilson Sonsini Goodrich & Rosati
- Orrick
- 10X Genomics
- GNS Healthcare
- Gerson Lehman Group
- Coatue Management
- Other corporate relationships
 - Northrop Grumman
 - Genentech

- Johnson and Johnson
- Optum
- Shares or Ownership
 - NuMedii (major)
 - Personalis (major)
 - Apple
 - Facebook
 - Alphabet (Google)
 - Microsoft
 - Amazon
 - Snap
 - 10x Genomics
 - Illumina
 - Nuna Health
 - Assay Depot (Scientist.com)
 - Vet24seven
 - Regeneron
 - Sanofi
 - Royalty Pharma
 - AstraZeneca
 - Moderna
 - Biogen
 - Paraxel
 - Sutro

- Speakers' bureau
 - None
- Companies started by students
 - Carmenta
 - Serendipity
 - Stimulomics
 - NunaHealth
 - Praedicat
 - MyTime
 - Flipora
 - Tumbl.in
 - Polyglot
 - lota Health
 - Ongevity Health

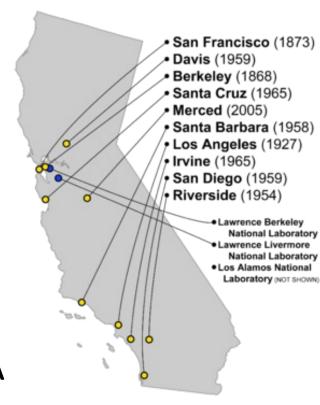
University of California

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

UC Health

- 20 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





Working with UCSF and UC-wide data

patient populations



- Access to deidentified, limited, and identified structured and unstructured clinical data
- Access to additional datasets that are or can be linked to clinical data (e.g., imaging, 'omics, waveforms)

- Access to deidentified structured clinical data from a larger patient population ()
- Access to additional data sets possible (---)
 but requires additional work and time (e.g.,
 recruitment and coordination of PIs at multiple
 sites, increased contractual complexity)

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 JNIVERSITY OF CALIFORNIA HEALTH



- Central database built using open-source OMOP as a data backend
 - First EHR installation was January 2012
 - Structured data from 2012 to the present day
 - 8.7 million patients with "modern" data
 - 378M encounters, 1.0B procedures, 1.3B+ med orders, 44M device uses,
 1.4M providers, 1.1B diagnosis codes, 5.2B+ lab tests and vital signs
 - "From Tylenol to CAR-T cells..."
 - Merged with California state data, pathology and radiology text elements, CA death index
 - Claims data from our self-funded plans now included
 - Continually harmonizing elements
- Safe, respectful, regulated, responsible use of clinical data

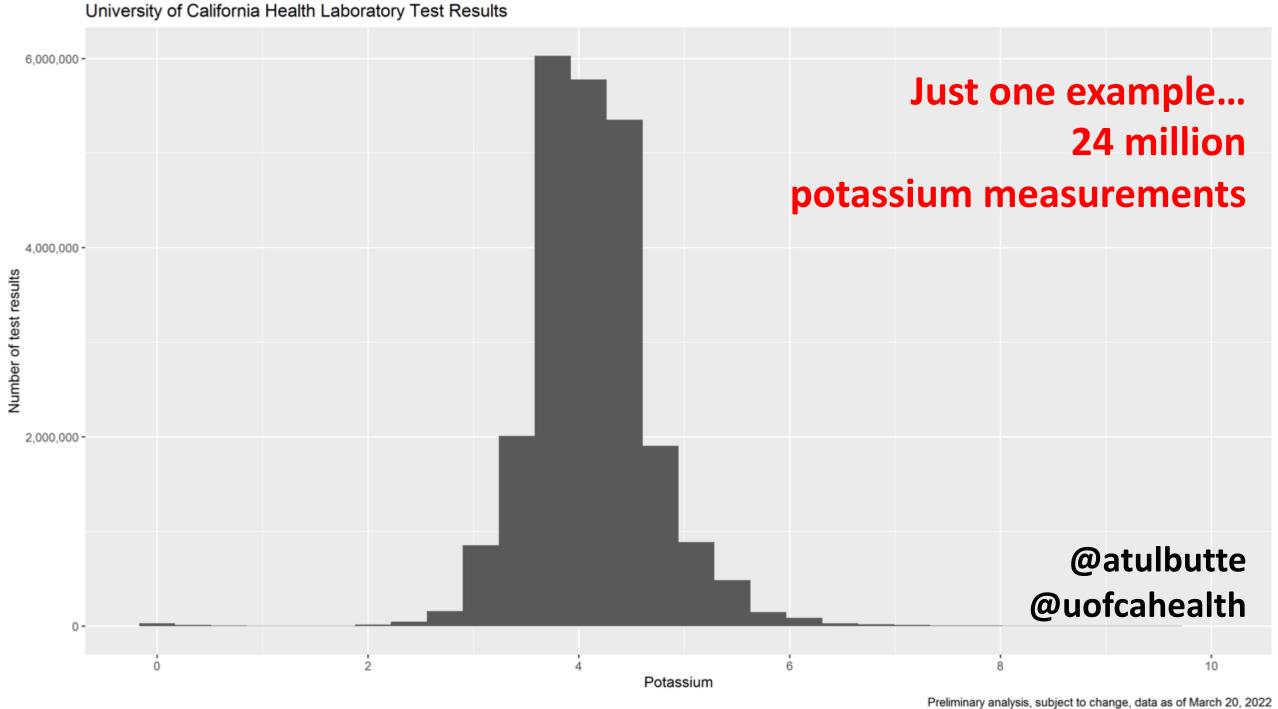




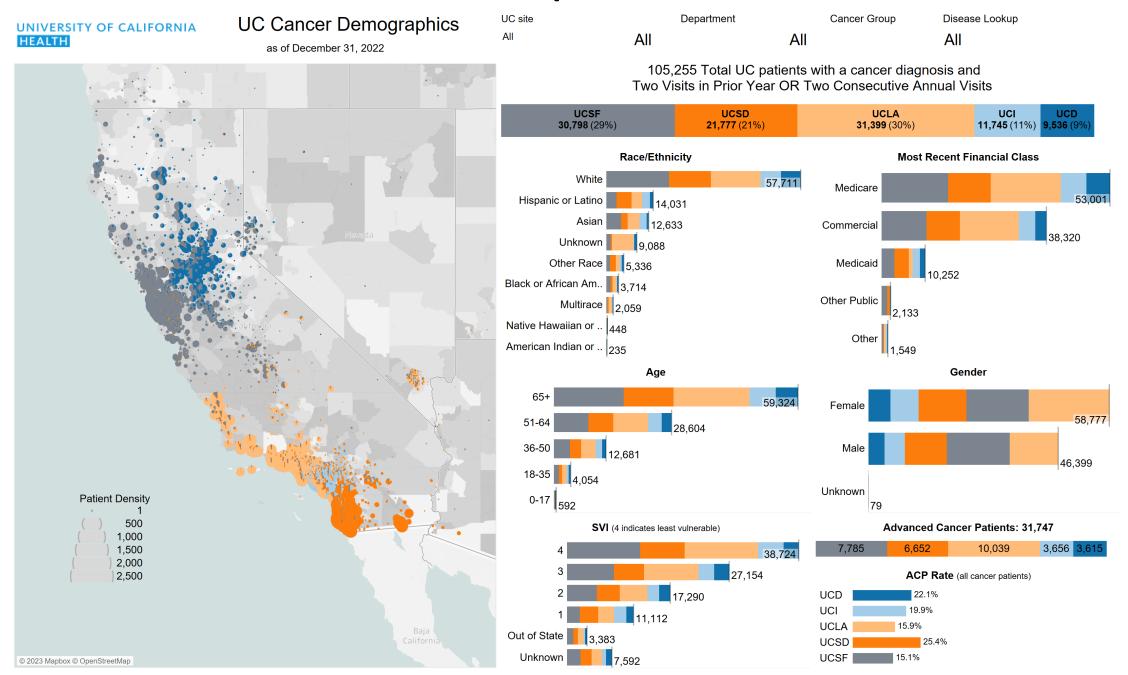




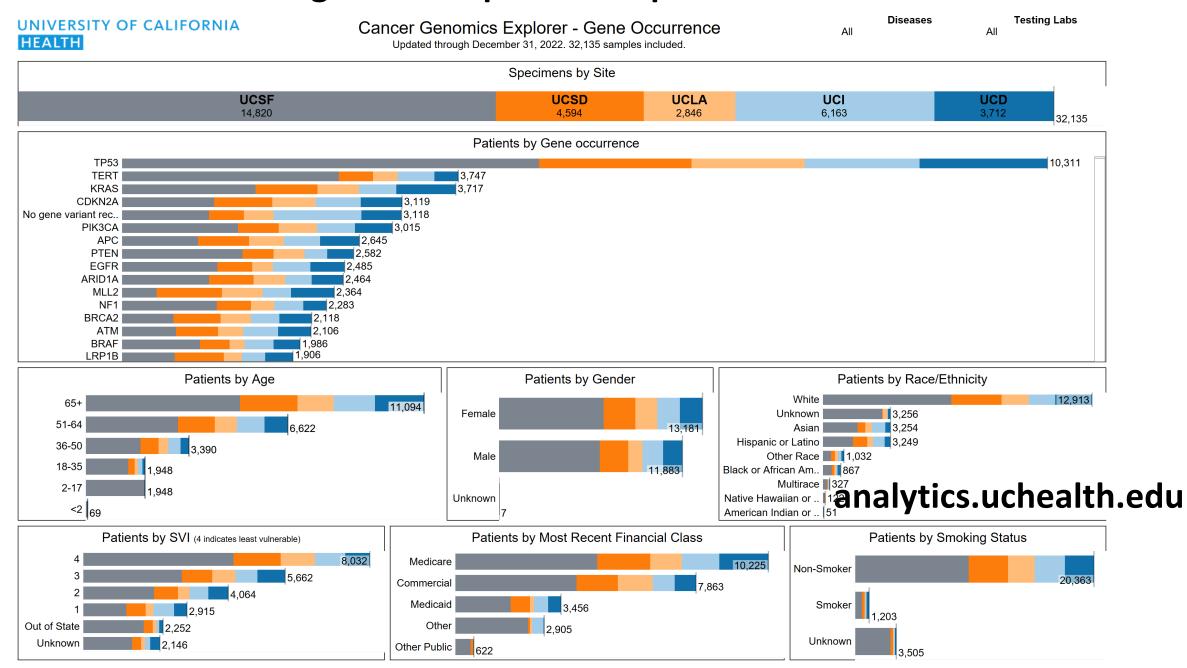
UC San Diego Health



100,000+ active cancer patients across UC Health



32+ thousand cancer genomic reports incorporated into the same database



View by cancer or by gene mutation

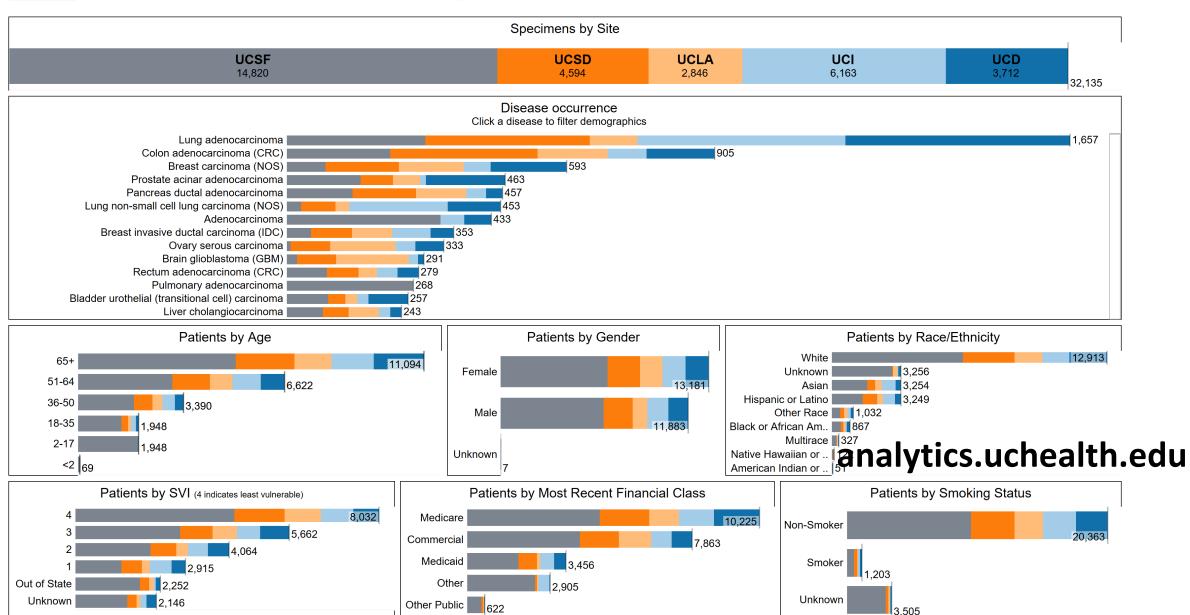
UNIVERSITY OF CALIFORNIA HEALTH

Cancer Genomics Explorer - Disease Occurrence

Genes

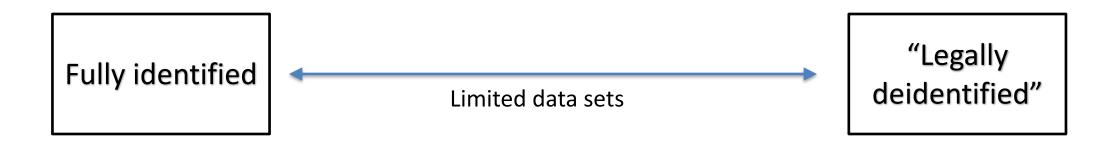
Testing Labs

Updated through December 31, 2022. 32,135 samples included.



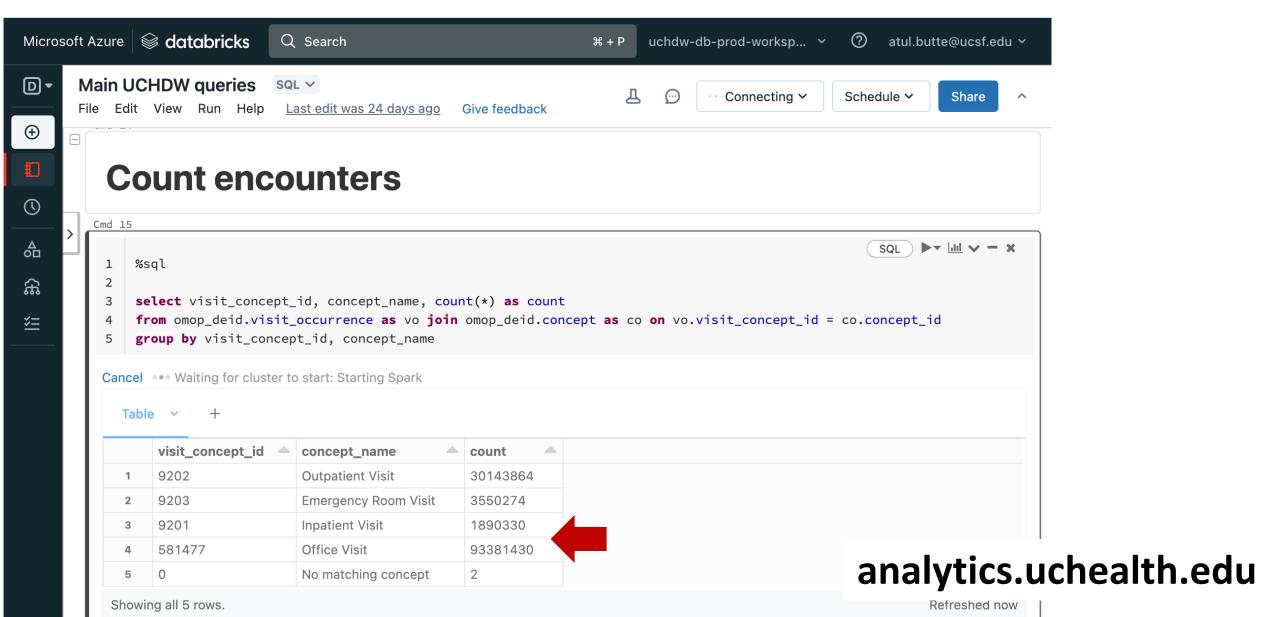
For example, only the cancers with an NF1 mutation **Testing Labs** Genes UNIVERSITY OF CALIFORNIA Cancer Genomics Explorer - Disease Occurrence NF1 HEALTH Updated through January 31, 2023. 32,797 samples included. Specimens by Site **UCSF** UCD **UCSD** UCLA UCI 383 1,186 419 274 298 2.560 Disease occurrence Click a disease to filter demographics Lung adenocarcinoma Lung non-small cell lung carcinoma (NOS) Breast carcinoma (NOS) Colon adenocarcinoma (CRC) Brain glioblastoma (GBM) Metastatic melanoma 45 Skin melanoma Melanoma Ovary serous carcinoma | Unknown primary melanoma Glioblastoma, likely IDH-wildtype, WHO grade IV Breast invasive ductal carcinoma (IDC) Lung squamous cell carcinoma (SCC) Unknown primary carcinoma (NOS) Patients by Gender Patients by Race/Ethnicity Patients by Age White 1.267 1,145 Unknown 51-64 Female Hispanic or Latino 36-50 Other Race 111111 18-35 Black or African Am.. 82 Native Hawaiian or ... 10 analytics.uchealth.edu 2-17 Male <2 1 American Indian or .. 6 Patients by SVI (4 indicates least vulnerable) Patients by Most Recent Financial Class Patients by Smoking Status 754 Medicare Non-Smoker Commercial Medicaid Smoker Out of State Unknown Unknown Other Public 59

How does research access work across UC Health?



- Researchers should first write and optimize OMOP SQL queries locally, on their own campuses
- When ready to scale (and authorized), we spin up a virtual machine for the researcher, populated with common tools
 - Electronically sign a UC Health data use agreement
 - R, Tableau, Jupyter Notebooks, Julia, SQL, Windows or Linux available
- Upload your scripts and run, but cannot download data
- Safe, respectful, regulated research use of clinical data

Safe, Respectful access to Deidentified Data now available through Cloud-based Databricks (single sign on with UC Health campus credentials)



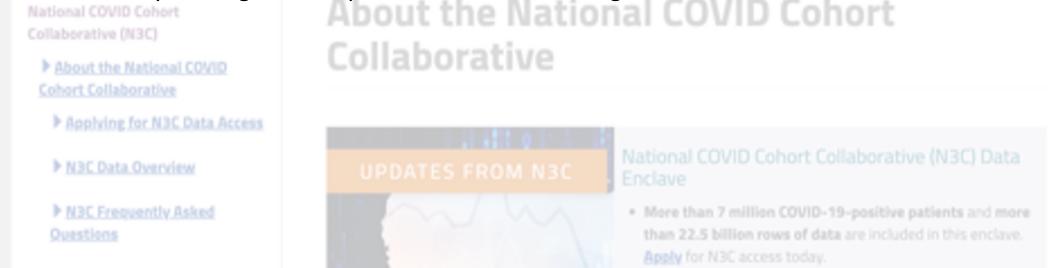
NCATS/NIH National COVID Cohort Collaborative (N3C)

- How was NCATS able to motivate 100+ health systems (especially their CTSA programs) to sign and transfer deidentified EHR data around SARS-CoV-2 tested individuals?
 - Common national IRB for reliance, for submission

About Translation

About NCATS

- Local IRB for scientific questioning, NIH Data Access Committee for access
- 463 publicly listed data projects underway (covid.cd2h.org/projects)
- 53 publications to date using this data (covid.cd2h.org/dashboard/index.jsp?publications)
- NCI still can't get this to happen for their designated cancer centers?
 - Do sticks work better than carrots for data sharing?
 - Are we continually waiting for data perfection before sharing?











N3C Dashboards

Reliable High-Velocity COVID-19 Insights Brought to you by the N3C.

Explore Dashboards

Publications

The N3C Data Enclave represents one of the largest secure collections of harmonized clinical health data in the United States.



Persons: 18.2 million

COVID+ Cases: 7,152,385

of Rows: 22.9 billion

Clinical Observations: 2.1 billion

Lab Results: 10.9 billion

Medication Records: 3.6 billion

Procedures: .9 billion

Visits: 1.2 billion

Explore Our Dashboards







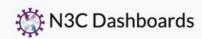


Medications









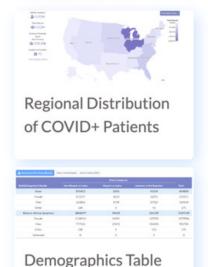


Dashboards > Institutions Contributing Data

Institutions Contributing Data

The N3C allows medical sites within the United States to securely transfer anonymized data into the Enclave. The average interval for data transfer from our partners is once a week. To explore the geographic coverage of our current partners, please see the map below.

Related Dashboards

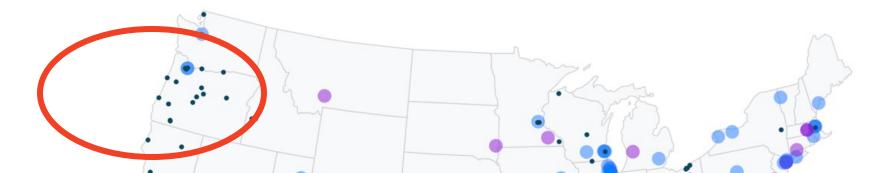


for IRB Submission

N3C Contributing Sites | 9



- Data Available
- Data transfer signed, pending availability
- OCHIN contributing site





Ongoing challenges and opportunities

- Health data interoperability/interchange is solvable, especially if there's a business reason for it
- Structured data elements from unstructured data is still a challenge
 - Expensive human curators now, artificial intelligence soon enough
- Mapping patient care trajectories to actual protocols (and deviations) is still hard
- Just because you have all the clinical and molecular data in one place does not mean people use it
 - Need better (easier) tools on top of these datasets, ideally within other tools being used (EHR)
 - Need to integrate the use of computational/data tools into the molecular cancer workflow
- Still hard to inspire cross-campus clinical trials: incentives not enough to push against all the resistance?
- We essentially have zero cancer precision medicine tools for patients
- Can we learn from the NIH/NCATS N3C experience?

UCSF Health

RICK LARSEN
OKSANA GOLOGORSKAYA
NELSON LEE
SANA SWEIS
THOMAS HURAY

UC San Diego Health

JENNIFER HOLLAND
DOUGLAS MACLEOD
PETER RYAN
HIRAM CARDOZA
ANDY LUCAS
Priscilla Jayaprakash

CALVIN FONG

UCI Health

DAVID MERRILL NORA LEWIN KATHY PICKELL LEANIE MAYOR MELODY HILL NEAKTISIA LEE

UNIVERSITY OF CALIFORNIA HEALTH

CARRIE BYINGTON

ATUL BUTTE

CORA HAN

PAGAN MORRIS

MONTE RATZLAFF

MIKE KILPATRICK

LISA DAHM

ANDENET EMIRU

IENNIFER BENBOW

EMRICA AGOSSA

AYAN PATEL

AIDEN BARIN

CHAYA MOHN

RAY PABLO

TIM HAYES

DAVID GONZALEZ

ROB FOLLETT

TEJU YARDI

Nadya Balabanova

UCLA Health

ALBERT DUNTUGAN
ANDREW WEAVER

YAEL BERKOVICH

JAY SHAH

VAJRA KASTURI

PALLAVI MYNAMPATI

SAAJID FAZIL

BILL LAZARUS

EDGAR TIJERINO

BILL CINNATER

JOSH PELINO

UNIVERSITY
OF
CALIFORNIA
HEALTH

Center for Data-driven Insights & Innovation

The CIO Team

Thank

you!

ELLEN POLLACK (UCLA)
JOSH GLANDORF (UCSD)

SCOTT JOSLYN (UCI)

JOE BENGFORT (UCSF)

ASHISH ATREJA (UCD)

Tom Andriola (UCI)

HEALTH

Kent Anderson
Puneet Gill
Hemanth Tatiparthi
Supraja Radhakrishnan
Jodi Nygaard
Jeffrey Sterett
Ralph Perrin

THOMAS AMI

Support

- University of California, San Francisco
- Priscilla Chan and Mark Zuckerberg
- Barbara and Gerson Bakar Foundation
- NIH: NIAID, NLM, NIGMS, NCI, NHLBI, OD; NIDDK, NHGRI, NIA, NCATS, NICHD
- Food and Drug Administration
- California Governor's Office of Planning and Research
- Howard Hughes Medical Institute, California Institute for Regenerative Medicine
- March of Dimes, Juvenile Diabetes Research Foundation
- Hewlett Packard, L'Oreal, Progenity, Genentech, Janssen
- Intervalien Foundation, Leon Lowenstein Foundation, Scleroderma Research Foundation, Clayville Research Fund, PhRMA Foundation, Stanford Cancer Center, Bio-X, SPARK
- Tarangini Deshpande
- Kimayani Butte
- Carrie Byington, Talmadge King, Mark Laret
- Jack Stobo, Sam Hawgood, Keith Yamamoto
- Isaac Kohane

Admin and Tech Staff

- Boris Oskotsky
- Andrew Jan
- Mounira Kenaani
- Andrew White
- Pam Allarde
- Amber Nolan