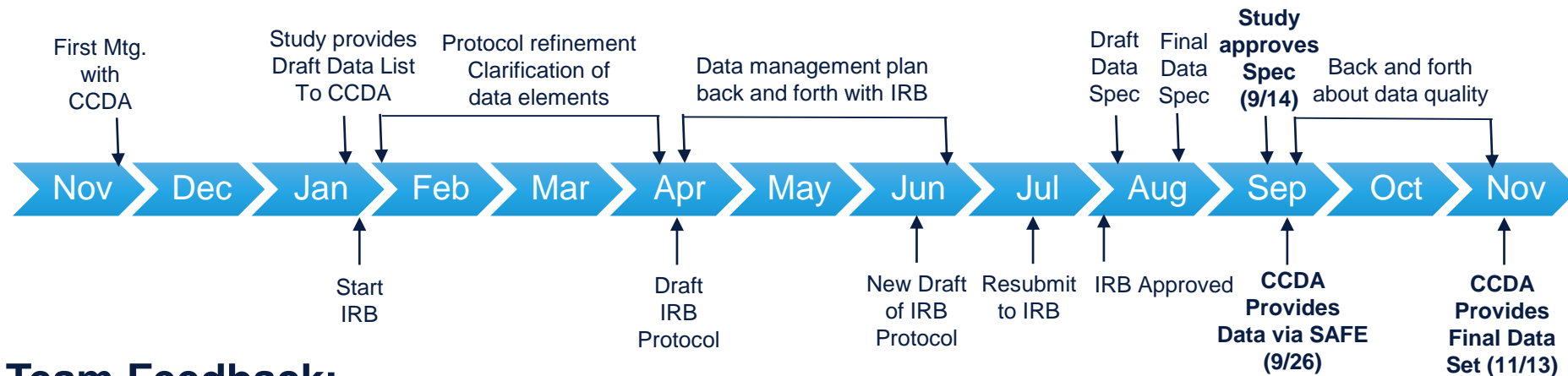# Precision Medicine
# *Analytics Platform*

*Paul Nagy, PhD, FSIIM*
*Associate Professor of Radiology*
*Division of Health Science Informatics*
*Armstrong Institute for Quality and Patient Safety*
*Deputy Director of JHM Technology Innovation Center*
*pnagy@jhu.edu*

JOHNS HOPKINS
APPLIED PHYSICS LABORATORY

JOHNS HOPKINS
UNIVERSITY

JOHNS HOPKINS
MEDICINE

On Average
4.5 mo. start-finish
20% on data extraction

# Case Study – Research Data Access

**Timeline & Feedback from CCDA Customer**

First Mtg. with CCDA

Study provides Draft Data List To CCDA

Protocol refinement Clarification of data elements

Data management plan back and forth with IRB

Draft Data Spec

Final Data Spec

**Study approves Spec (9/14)**

Back and forth about data quality

| Nov | Dec | Jan | Feb | Mar | Apr | May | Jun | Jul | Aug | Sep | Oct | Nov |

Start IRB

Draft IRB Protocol

New Draft of IRB Protocol

Resubmit to IRB

IRB Approved

**CCDA Provides Data via SAFE (9/26)**
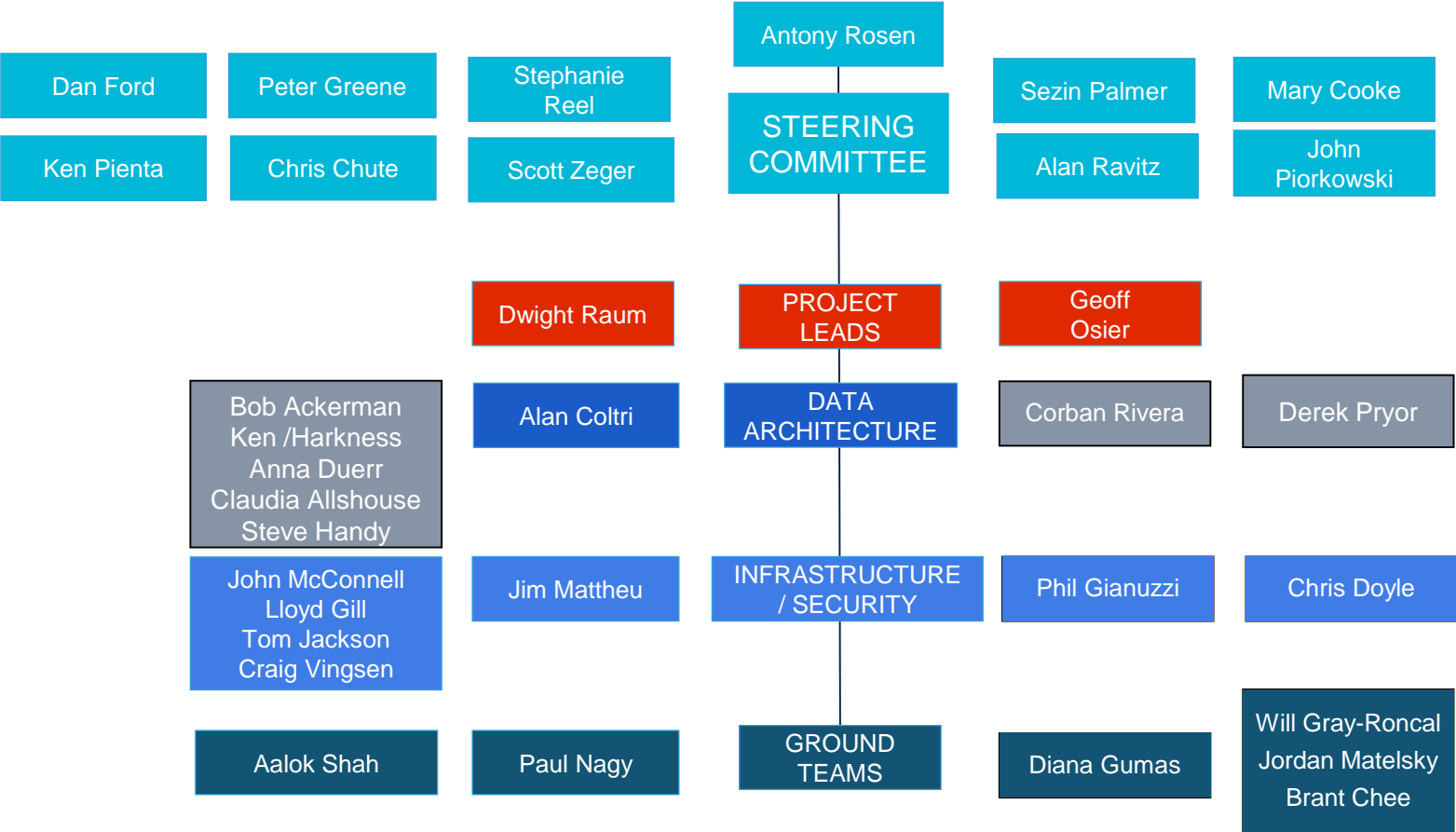
**CCDA Provides Final Data Set (11/13)**

**Study Team Feedback:**

1. **No** comprehensive **overview of the entire process**. We stumbled our way through.

2. Not much guidance on **how to complete the various IRB forms**. A **standard template** would help.

3. **Review (of data management plan for IRB**) is a rate-limiting step because it is only done by one person.

4. Researchers have to do a lot of work trying to figure out exactly what data elements they need and then provide the CDDA. A **meta-thesaurus would be useful**.

5. Researchers do not know **which data can be obtained easily** and those that cannot.

JOHNS HOPKINS **in**health

*in*Health is using revolutionary tools of measurement, data science, and connectivity to discover clinically-relevant and biologically-anchored subgroups at scale, and to deliver what we learn to impact the precision and value of health care

# School of Medicine and Applied Physics Laboratory



Antony Rosen

Dan Ford | Peter Greene | Stephanie Reel | STEERING COMMITTEE | Sezin Palmer | Mary Cooke
Ken Pienta | Chris Chute | Scott Zeger | | Alan Ravitz | John Piorkowski

Dwight Raum | PROJECT LEADS | Geoff Osier

Bob Ackerman
Ken /Harkness
Anna Duerr
Claudia Allshouse
Steve Handy | Alan Coltri | DATA ARCHITECTURE | Corban Rivera | Derek Pryor

John McConnell
Lloyd Gill
Tom Jackson
Craig Vingsen | Jim Mattheu | INFRASTRUCTURE / SECURITY | Phil Gianuzzi | Chris Doyle

Aalok Shah | Paul Nagy | GROUND TEAMS | Diana Gumas | Will Gray-Roncal
Jordan Matelsky
Brant Chee

Large, multidisciplinary, cross-functional team

# Changing Face of Medicine: Centers of Excellence

# Tenets of PMAP

1. Researchers need better access to clinical data.

2. Researchers need environments that ensure data security protecting patient information recognizing full de-identification is difficult

3. Researchers need an environment that is built for machine learning and data science to enable discovery.

4. Combined access to very different data types: EMR, medical imaging, genomics, and physiological monitoring data.

5. Clinical researchers need to bring new discoveries into clinical care.

# Platform Components

## Data Platform

- Confidentiality
- Integrity
- Availability
- Authentication
- Authorization
- Accounting

## Tooling

- Patient matching
- Data Catalog
- Cohort Discovery
- Honest Broker
- Annotation
- Preprocessing

## Research Projections

- SQL Server
- Cohort Dashboard
- Jupyter Notebooks
- Docker
- Compute

# Make JHM Easy

- Risk Tiers – Bioethical Framework
  - Proposed by IRB and Data Trust
  - Created to accommodate data science in PMAP
  - Leverages secure analysis environment (SAFE)

- Tier A Proposals
  - Approved as a class by Data Trust
  - IRB review streamlined
  - Investigators reference data categories rather than data elements

- Positive Investigator Impact
  - Simplified data specification (save hours)
  - Faster review time (save months)

# NLP

- Current research being done with Prostate Cancer CoE
- \> 95% accurately extracts Gleason scores and anatomical references from notes
- Post-processing on data allows
  - Large-scale inference
  - Error detection and other data quality downstream analytics
- Beginning to evaluate using this for other CoEs

# Medical Images VNA Access

- PMAP Data Commons only stores the imaging metadata (DICOM)
  - Patient, study, series, image
- Users can query DICOM with Hive's SQL-like language
- Users (with appropriate permissions) request images to be fetched for their Projections
- Images pulled from VNA
- Deep learning GPU Compute



VNA DICOM Data

# Genomics on PMAP



**On platform vs. Federation**
- Caching strategy: cache gene and allele level features on the platform. Federate access to raw sequence
- Cached formats: VCF, sample metadata, Gene expression matrices, genome references
- Federated: fastq, bam

**Public Annotation Sources**
- ClinVar, OMIM, GeneOntology

**Functionality Targets**

Cohort Discovery
> Semantic search that allows the user to discover patients based on variants/ variant sets or genomic alterations (indels, CNV).

Cohort Extraction
> A method for displaying/exporting variants and associated metadata for that group of patients returned by a search.

Data Explorer
> A genome browser that allows the user to drill down and display the aligned sequence data and associated quality metrics in the chromosomal location defined by the user (around variants).

Data Commons
> Genomics derived datasets linked by a patient identifier.
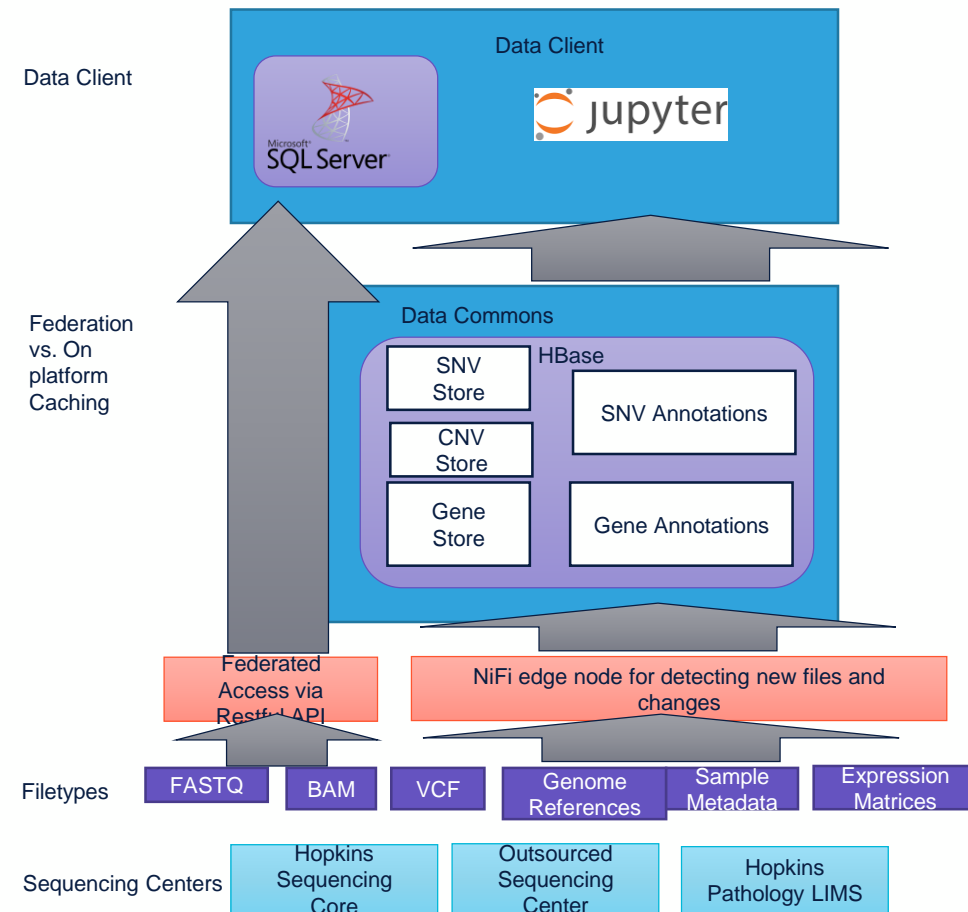
Variant Annotation
> The ability to add characterizations to a variant, according to ACMG scales or other well-defined methods

Variant Filtering
> The ability to filter a returned list of variants by criteria or metadata associated with the sequence.

Container Environment for Analysis Scaling
> The ability to configure the container environment based upon the study needs

# Physiological Monitoring

- Hadoop environments are well suited for large streaming data real time analysis environments.

- Twitter Firehouse:
  - 500 Million tweets a day

- PhysioCloud research group.

# Overall Layout

## Public Resources

### Precision Medicine Introduction

What are inHealth, PMCOEs, and the PMAP, and why you would want to use these resources

News and Events: Relevant Conferences, Programs, Newly published resources / code examples

### Researcher Resources

Becoming a PMCOE

- PMCOE in a Box Templates
- e.g. Value Matrix, Cohort in a Box

Templates for grants/IRB/Data Trust applications

Research Lifecycle (or link to ICTR resource)

Research Community Resources (e.g. browsing data scientist / faculty profiles, contact ICTR/CCDA/TIC)

### PMAP Education

Documentation, Videos on using the technologies
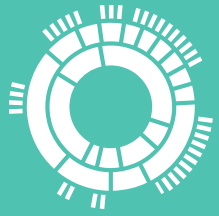
Gallery of code / notebooks

## Community Only (JHED Login)

### PMAP Administration

- Sign up
- Account Management (Tiering, Usage reporting)
- Payment
- Team Management / Role Assignment
- Define databases to ingest into PMAP
- Request Data Projections
- Help / Bug Reporting / Open Issues

### PMAP Tools

- Data Catalog
- Cohort Discovery
- Cohort Dashboard
- Crunchr
    - Access to compute resources, standard containers, shared notebooks, version control through Crunchr command line
- Data Annotation Tools
    - NLP, Imaging, Genomics annotation tools
- Projection Creator (select users only, e.g. CCDA)

# CAMP
## Center of Excellence Analytics in Medicine Program

Opportunity for PMAP team to work directly with the JHMI research teams

# CAMP - CoE Analytics in Medicine Program

## 12 week course designed to introduce researchers to PMAP

Opportunity for PMAP team to work directly with the JHMI research teams

Provide an overview of PMAP's infrastructure, data management, data science capabilities, and clinical applications

Further understanding of what tools will be most useful

Work together on research proposals

*Non-PMCoE teams are charged*

# CAMP - CoE Analytics in Medicine Program

- ✔ Data Needs
- ✔ PMAP Overview
- ✔ Working With IRB / Data Trust
- ✔ Working With CCDA
- ✔ Epic Data Sources
- ✔ Machine Learning For Imaging

- ✔ Derived Features For Genomics
- ✔ Physio/Wearables Data Use
- ✔ Machine Learning For NLP
- ✔ Overview of Partner Institutions
- ✔ Cohort Dashboard
- ✔ Integrate Research Into Clinic

# CAMP Team



CAMP
2018 FACULTY

**Paul Nagy, PhD**
Deputy Director, TIC
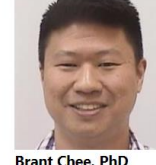*CAMP Director*

**Diana Gumas**
Senior IT Director, ICTR
*CAMP Director*

**Ken Pienta, MD**
Professor, Urology Research
*CAMP PMCOE Lead*

**Alan Coltri**
Director, Data Arch +
Integration, IT@JH CIO

**Brant Chee, PhD**
APL

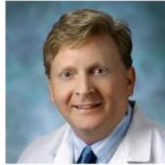**Dwight Raum,**
Chief Technology
Officer, Johns Hopkins

**Ken Harkness**
Project Lead, IT@JH
Clinical Systems Dev

**Patrick Ostendarp**
Product Development
Lead, TIC

**Corban Rivera, PhD**
APL
*CAMP APL Lead*

**Gregory Kirk, MD**
Professor of Epidemiology, SPH
*CAMP Research Lead*

**Aalok Shah**
IT Product Development Manager,
TIC
*CAMP Manager*

**Emily Marx**
Communications, TIC
*CAMP Manager*

**Alex Baras, MD**
Assistant Professor,
Pathology Informatics

**Caitlyn Bishop**
APL

**Jerry Prince, PhD**
Professor, Electrical +
Computer Engineering

**Lee Watkins**
IT Operating Unit
Director, CIDR

**Rai Winslow, PhD**
Director, Biomedical
Engineering

**Scott Zeger, PhD**
Director, inHealth
*CAMP Advisor*

**Antony Rosen, MD**
Vice Dean for Research, SOM
*CAMP Advisor*

**Mary Cooke**
Director, inHealth
*CAMP Advisor*

**Chris Chute, MD**
Chief Health Research Info Officer,
JHM
*CAMP Advisor*

**Benjamin Smith**
Application Coordinator,
IT@JH Epic

**Chris Doyle**
Product Development
Lead, TIC

**John Scott**
Senior Software
Engineer, TIC

**Manar Alhamdy**
Senior Software
Engineer, TIC

**Stephen Granite**
Senior IT Manager, ICM

**Bonnie Woods**
Senior IT Manager,
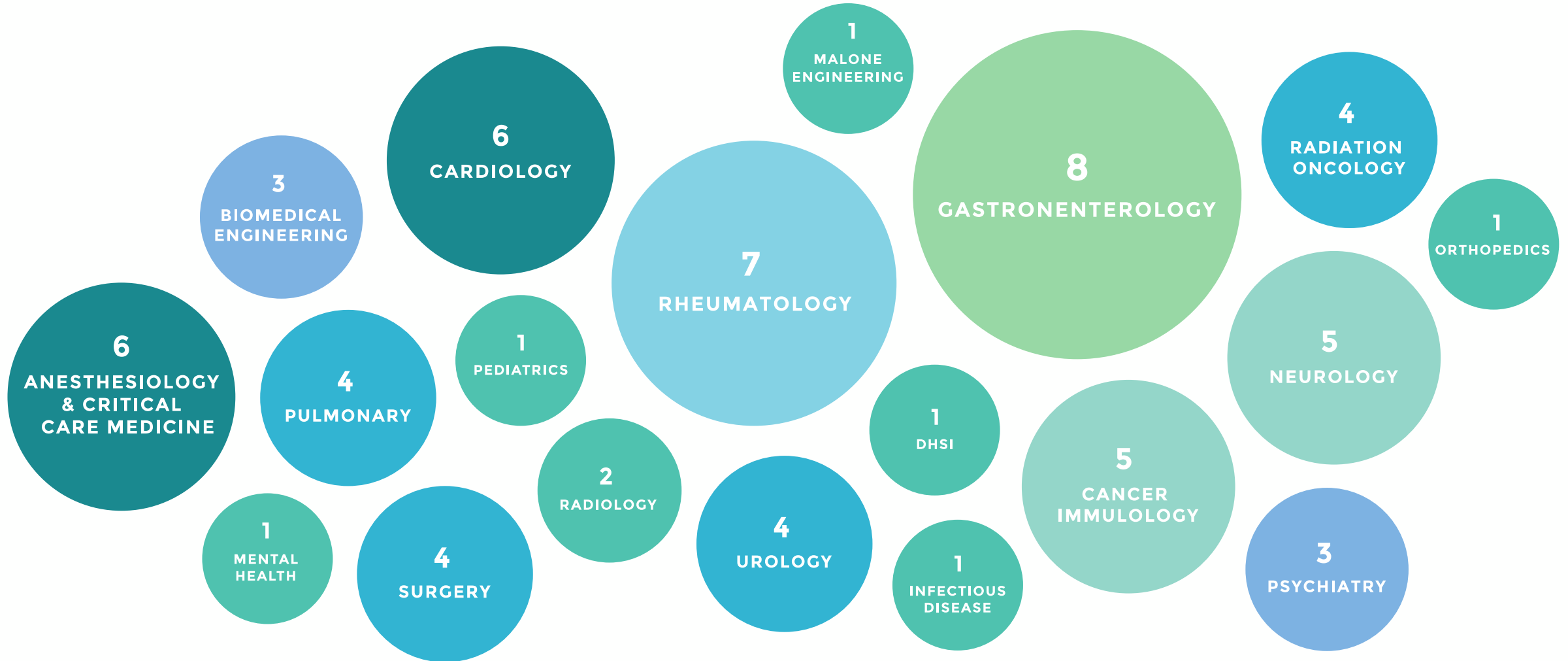ICTR

**David Li**
IT Director, IT@JH Epic

**Jordan Matelsky**
APL

**Masoud Rouhizadeh, MD**
Software Engineer, ICTR

**Steven Handy**
Senior Software
Engineer, TIC

# Department Distribution

- 1 MALONE ENGINEERING
- 6 CARDIOLOGY
- 3 BIOMEDICAL ENGINEERING
- 8 GASTRONENTEROLOGY
- 4 RADIATION ONCOLOGY
- 1 ORTHOPEDICS
- 7 RHEUMATOLOGY
- 6 ANESTHESIOLOGY & CRITICAL CARE MEDICINE
- 4 PULMONARY
- 1 PEDIATRICS
- 5 NEUROLOGY
- 1 DHSI
- 5 CANCER IMMULOLOGY
- 1 MENTAL HEALTH
- 4 SURGERY
- 2 RADIOLOGY
- 4 UROLOGY
- 1 INFECTIOUS DISEASE
- 3 PSYCHIATRY

# CAMP Makeup (N=70)

| | | |
|---|---|---|
| **12**<br>Full Professors | **13**<br>Associate Professors | **14**<br>Assistant Professors |
| **1**<br>Instructors | **10**<br>Fellows/Students | **20**<br>Research Staff |

# Jupyter: Open Data Science at Scale



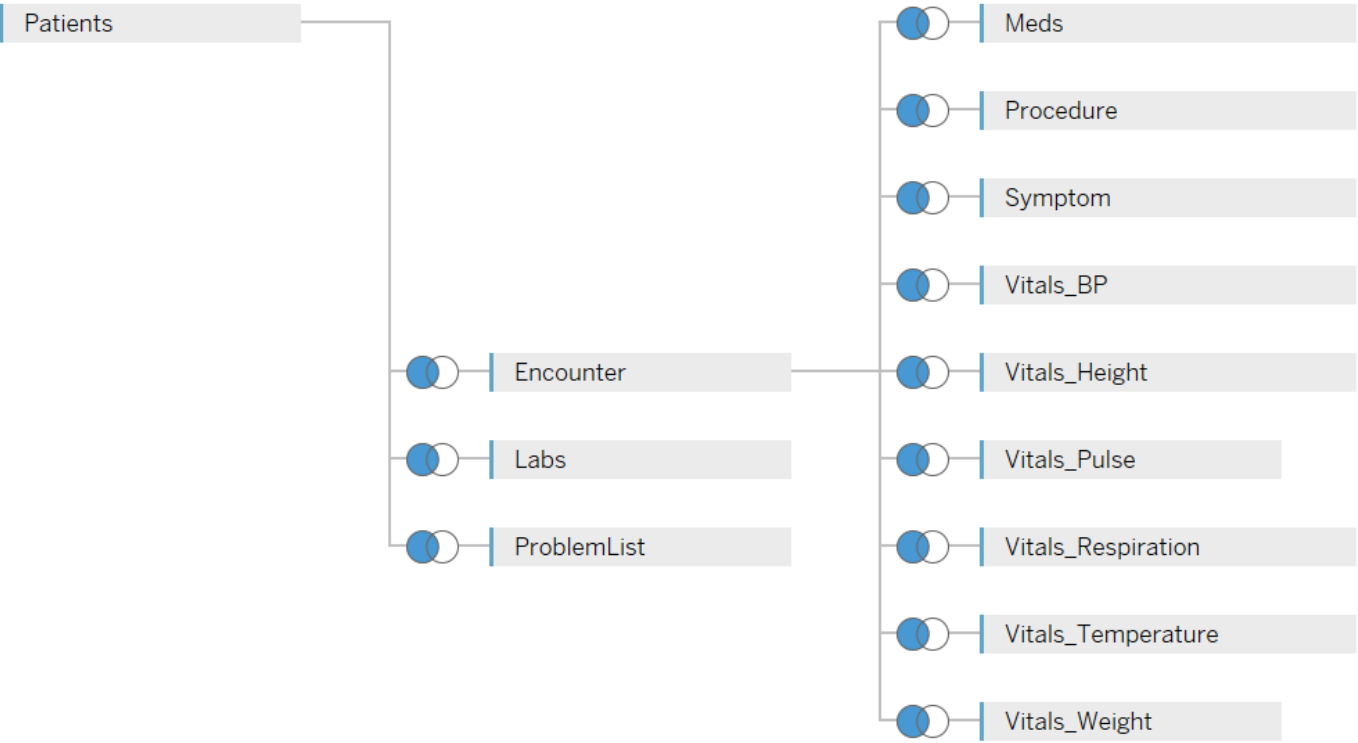**Rich web client** ➡ **Text & Math** ➡ **Code** ➡ **Results** ➡ **Share, reproduce**

# De-identified Epic Dataset



112 Million Data Elements
No note fields

| Table Name | # of Records | # of Columns | Size |
|---|---|---|---|
| Encounter | 690,183 | 4 | 36.57 MB |
| Labs | 3,486,911 | 12 | 403.41 MB |
| Meds | 5,926,733 | 9 | 1.86 GB |
| Patients | 60,676 | 5 | 7.18 MB |
| ProblemList | 115,162 | 4 | 16.99 MB |
| Procedure | 6,449 | 6 | 0.727 MB |
| Symptom | 28,056 | 5 | 1.711 MB |
| Vitals_BP | 390,181 | 7 | 30.58 MB |
| Vitals_Height | 279,288 | 7 | 18.96 MB |
| Vitals_Pulse | 388,450 | 7 | 65.92 MB |
| Vitals_Respiration | 251,166 | 7 | 18.18 MB |
| Vitals_Temperature | 314,571 | 7 | 21.87 MB |
| Vitals_Weight | 352,553 | 7 | 24 MB |
| TOTAL | 12,290,379 | 87 | 2.55 GB |

# Platform Architecture



**DATA SOURCES**

Pathology

Epic

Radiology

.
.
.

**Includes local databases**

Reusable Pipelines

DISCOVERY

**DATA COMMONS**

Hosts all data coming into platform from various, disparate data sources

IRB Approvals

**RESEARCH ENVIRONMENT**

Secure user environment with aggregated data sets approved for user – allows advanced analytics and other tools to be applied to data

Validation/ Promotion Process

FUTURE DELIVERY

**TREATMENT PLATFORM**

Shared Algorithms, Analytics, Clinical Care Impact, Products

Commercialize

Feedback

# Decision Support



Patient test results

Expert Knowledge

Scientific Literature

Predictive Models

Clinical Data Warehouse

Shared decision making

Patient goals and preferences