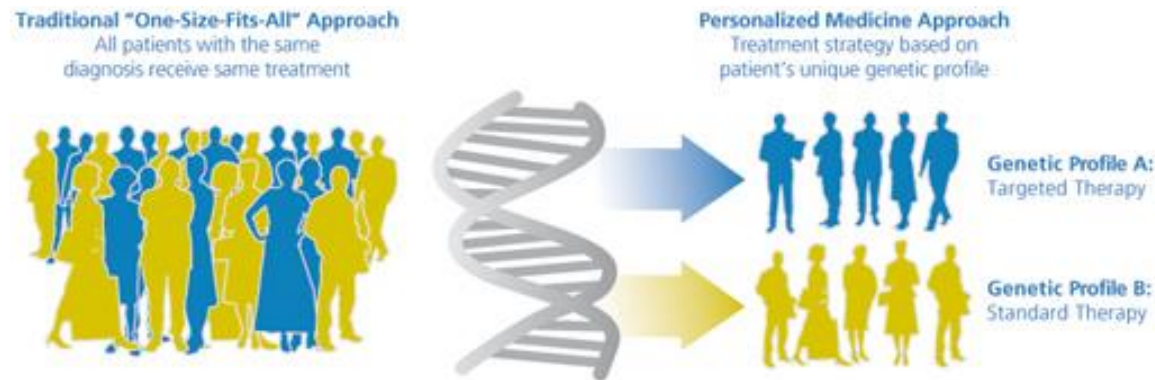# Instrumenting the Health Care Enterprise for Discovery in the Course of Clinical Care

*Shawn Murphy MD, Ph.D.*

*Chief Research Information Officer*
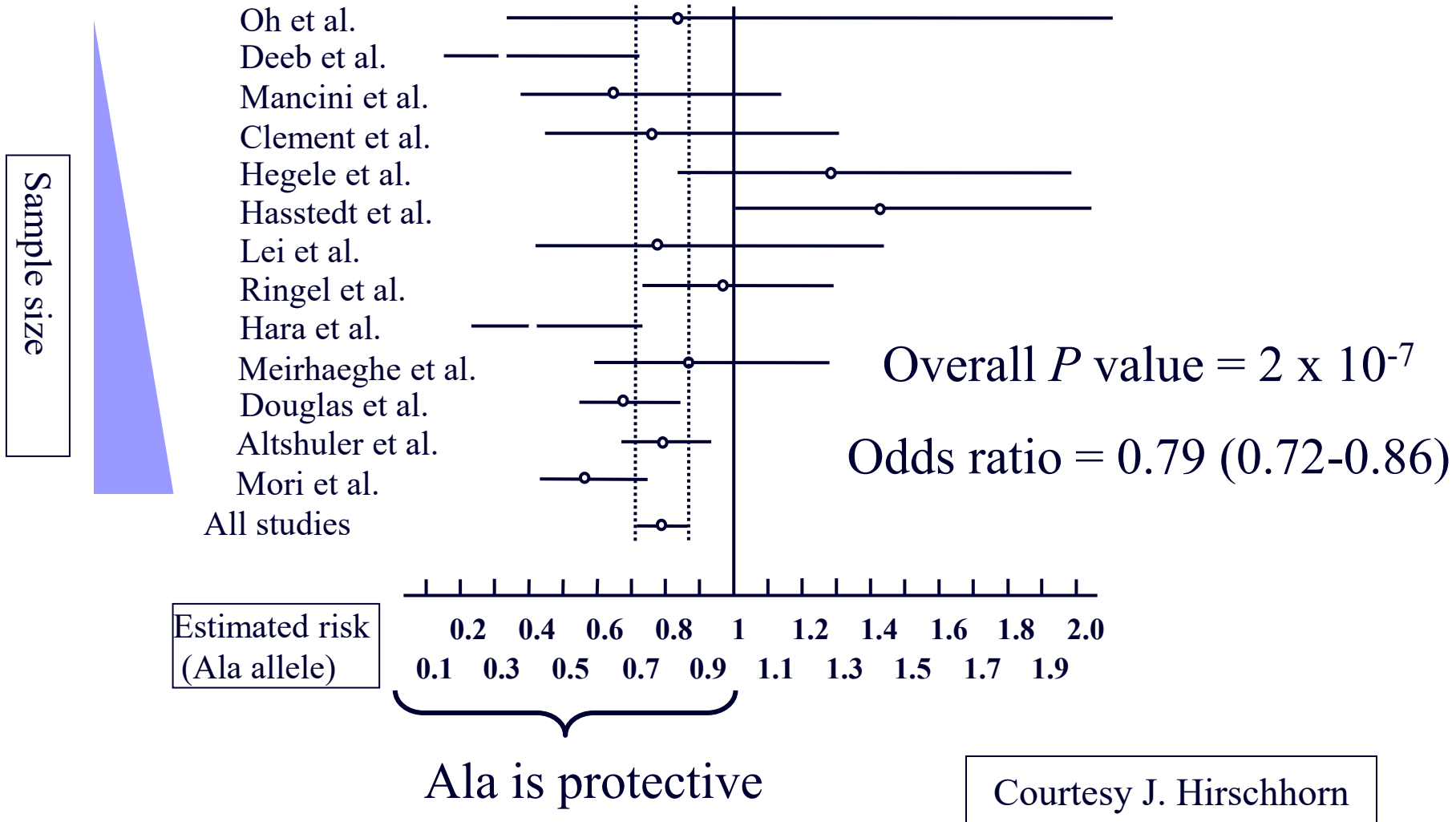
*Harvard Medical School / Mass General Brigham*

# Personalized Medicine and Genomic technology are critical to managing populations



Traditional "One-Size-Fits-All" Approach
All patients with the same diagnosis receive same treatment

Personalized Medicine Approach
Treatment strategy based on patient's unique genetic profile

Genetic Profile A:
Targeted Therapy

Genetic Profile B:
Standard Therapy

- Managing a population involves improving health outcomes of the group as a whole by identifying, monitoring and addressing health needs of individuals through:
  - Subpopulation stratification
  - Targeted, evidence-based treatment protocols
  - Predictive analytics

# Example: PPARγ Pro12Ala and Diabetes



Overall *P* value = $2 \times 10^{-7}$

Odds ratio = 0.79 (0.72-0.86)

Sample size

Estimated risk (Ala allele)

0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9  1  1.1  1.2  1.3  1.4  1.5  1.6  1.7  1.8  1.9  2.0

Ala is protective

Courtesy J. Hirschhorn

Oh et al.
Deeb et al.
Mancini et al.
Clement et al.
Hegele et al.
Hasstedt et al.
Lei et al.
Ringel et al.
Hara et al.
Meirhaeghe et al.
Douglas et al.
Altshuler et al.
Mori et al.
All studies

# High Throughput Methods for supporting Translational Research

- Set of patients is selected from medical record data in a high throughput fashion

- Investigators explore phenotypes of these patients using Machine Learning tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Digital medicine is delivered into clinical care through Digital Twin

# Data problems that make working with Electronic Healthcare Data to conduct research difficult

- 1) There are significant risks of a data breach which will result in very large fines and loss of confidence in the hospitals where the breach occurred.

- 2) The data are not collected for research purposes, and therefore the data can be poorly structured with significant omissions, biases, and inaccuracies.

# Research Patient Data Registry (RPDR) at Mass General Brigham to find patient cohorts and distribute data

## 1) Queries for aggregate patient numbers

- Warehouse of in & outpatient clinical data
- 6.7 million Mass General Brigham patients
- 2.6 billion diagnoses, medications, genomics, procedures, laboratories, & physical findings coupled to demographic & visit data
- Authorized use by faculty status
- Clinicians can construct complex queries
- Queries cannot identify individuals, internally can produce identifiers for (2)

**Query construction in web tool**



De-identified Data Warehouse

Z731984X
Z74902XX
...
...

Encrypted identifiers

## 2) Returns detailed patient data

- Start with list of specific patients, usually from (1)
- Authorized use by IRB Protocol
- Returns contact and PCP information, demographics, providers, visits, diagnoses, medications, procedures, laboratories, microbiology, reports (discharge, LMR, operative, radiology, pathology, cardiology, pulmonary, endoscopy), and images into a Microsoft Access database and text files.

0000004
2185793
...
...

## OR

0000004
2185793
...
...

Real identifiers

# FINDING PATIENTS



Query items

Person who is using tool

Query construction

Results - broken down by number distinct of patients

# Theory of Kimball translated to Healthcare Data

## Star schema

**Concept DIMENSION**

concept_key
concept_text
search_hierarchy

**Patient-Concept FACTS**

patient_key
concept_key
start_date
end_date
practitioner_key
encounter_key
value_type
numeric_value
textual_value
abnormal_flag

**Encounter DIMENSION**

encounter_key
encounter_date
hospital_of_service

**Patient DIMENSION**

patient_key
patient_id (encrypted)
sex
age
birth_date
race
deceased
ZIP

**Pract . DIMENSION**

practitioner_key
name
service

## Binary Tree

start
search

.22

250

6.7

.04

2600 million

# RPDR Detailed Data Request Wizard -- Web Page Dialog

## RPDR DETAILED DATA REQUEST WIZARD

Using IRB#mgh-demo-1 (found in the RPDR Identified database) to obtain data from the RPDR

You are logged in as Murphy, Shawn N. in workgroup Shawn Murphy, MD

### Select protocol number(s)

Partners IRB (required): `mgh-demo-1`

Title: RPDR protocol - Demonstration IRB number for Dr. Murphy
Status: Active

Newton Wellesley Hospital IRB: `NWH Demo 1`

Title: test
Status: Active

Spaulding Rehabilitation Hospital IRB: `                    `

Options for returned set of patients:

☐ Create a static set of patients from this query that can be used in other RPDR queries

☑ Rerun the base query shown above to obtain a fresh set of patients

| Help | < Back | STEP 3 | Next > | Cancel |

**RPDR Detailed Data Request Wizard -- Web Page Dialog**

## RPDR DETAILED DATA REQUEST WIZARD

Using IRB#mgh-demo-1 (found in the RPDR Identified database) to obtain data from the RPDR

You are logged in as Murphy, Shawn N. in workgroup Shawn Murphy, MD

### Please select if you would like a HIPAA-defined (deidentified) limited data set or an identified data set

[ What's a limited data set? ]

○ **Limited Data Set**
  - The files that result from this request will be available in a protected file share with no special encryption.

⦿ **Identified Data Set**
  - The text files that result from this request will be encrypted and the Microsoft Access file will be password protected. In order to access the data, a password will be provided.

| Help | < Back | STEP 8 | Next > | Cancel |

**RPDR Detailed Data Request Wizard -- Web Page Dialog**

## RPDR DETAILED DATA REQUEST WIZARD
Using IRB#mgh-demo-1 (found in the RPDR Identified database) to obtain data from the RPDR
You are logged in as Murphy, Shawn N. in workgroup Shawn Murphy, MD

### Select the types of data that should be returned from the RPDR
### Only data allowed by your protocol should be chosen
(Identified data sets will always return a set of identified patient medical numbers)

**Detail Data Items**

- ☐ 📁 Demographic Data
- ☐ 📁 Identifying Patient Information - not available for Limited Data Sets
- 📁 LMR (Longitudinal Medical Record)
- 📁 Medications, Diagnoses and Procedures
- 📁 Medications, Diagnoses and Procedures from Billing Data - only visits where query criteria occur all in the same visit
- 📁 Patient Clinical Reports- not available for Limited Data Sets
  - ☐ Cardiology Reports
  - ☐ Discharge Summaries
  - ☐ Endoscopy Reports
  - ☑ Microbiology Data
  - ☑ Operative Notes
  - ☐ Pathology Reports
  - ☐ Pulmonary Reports
  - ☐ Radiology Reports
  - ☐ Transfusion Data, Blood Bank Data

[ Help ]  [ < Back ]  **STEP 9**  [ Next > ]  [ Cancel ]

# Detailed data is gathered for request and distributed



Output files placed in special directory

Data is gathered from RPDR and other MGB sources

Files include Small Database

# One year's usage of RPDR

- **4526 registered users, 1113 new in just 2019**

- **834 teams/year gathering data for research studies**

- **4472 detailed patient data sets returned to these teams in 2019, containing data of 24.7 million patient records.**

- **From a survey of 153 teams**
  - **Importance of the data received from the RPDR was evaluated in relation to the study it was supporting.**
  - **Calculated over 4 years (FY15-FY19) the total agreement amounts were $2.27 Billion, making per year consumption critically dependent on RPDR $244 Million.**



**Usefulness of Detailed Data**
*106 Total Responses*

Not Useful 15%
Critical 43%
Useful 42%

# *Rapid investigation of QTc prolongation*

■ *FDA warning 2011 for Celexa*

**Safety Announcement:**
**[8-24-2011] "should no longer be used at doses greater than 40 mg per day** because it can cause abnormal changes in the electrical activity of the heart."

■ *But, did NOT include Lexapro (which is active ingredient of Celexa [s-enantiomer])*

■ *Shown to be true with RPDR-derived data set with >38,000 EKGs obtained within 14 – 90 day window after medication initiated*

| Anti-depressant | Adjusted model† | |
|---|---|---|
| | prolongation | p-value |
| SSRI | | |
| Citalopram (Celexa) | 2.85 | 0.004 |
| Escitalopram (Lexapro) | 3.80 | < 0.001 |
| Fluoxetine (Prozac) | 1.44 | 0.150 |
| Paroxetine (Paxil) | 0.07 | 0.943 |
| Sertraline (Zoloft) | 0.87 | 0.383 |
| Other anti-depressants | | |
| Amitriptyline | 4.10 | < 0.001 |
| Bupropion | -2.15 | 0.032 |
| Duloxetine | 0.60 | 0.547 |
| Mirtazapine | -1.46 | 0.145 |
| Nortriptyline | 1.23 | 0.219 |
| Venlafaxine | 1.15 | 0.251 |
| **previously known prolonger** | | |
| Methadone | 5.32 | < 0.001 |

† Adjusted for age, gender, race, type of insurance, history of major depression, history of myocardial infarction and Charlson comorbidity score

Roy Perlis MD, MSc and team

# Relevant Cohorts of Patients are Gathered through RPDR and Detailed Data Obtained

- Medication use by individual patients over time
- Patient EKG QTc values at various time points

# Results: QTc interval and medication use



* Dose a significant predictor of QTc in fully adjusted linear models at α=0.05
† QTc at specified dose is significantly different from that at prior dose in fully adjusted linear models at α=0.05

Mean (SD) corrected QT (QTc) interval recorded on electrocardiogram 14–90 days after prescription of antidepressant or methadone, by drug dose

18

# High Throughput Methods for supporting Translational Research

- Set of patients is selected from medical record data in a high throughput fashion

- Investigators explore phenotypes of these patients using Machine Learning tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Digital medicine is delivered into clinical care through Digital Twin

**RPDR Evolved into international "Informatics for Integrating Biology and the Bedside (i2b2)" sponsored by the National Institutes of Health, what is it?**

- Software for explicitly organizing and transforming person-oriented clinical data to a way that is optimized for clinical genomics research
  - Allows integration of clinical data, trials data, and genotypic data
- A portable and extensible application framework
  - Software is built in a modular pattern that allows additions without disturbing core parts
  - Available as open source at https://www.i2b2.org

# I2b2 Community Software distributed as open source

# I2b2 Software adapts through new plugins

# Genotype Data



https://community.i2b2.org/wiki/display/IGD/Loading+Genomic+VCF+Files+into+i2b2

# Use NLP to extract the relevant features from the set of patient notes.

# LMM Enhanced interaction with Patient Representation

# Medical conditions supported by description in chart

# Data Integration in Big Data Commons

## Electronic Medical Record (EMR) Data

**RPDR**

Coded Data | Text Data (Notes/Reports)

| Demographics | Medications | Physician Notes |
| Diagnoses | Procedures | Imaging Reports |
| Lab Results | Visits | Pathology Reports |
| | | Surgery Notes |

## Additional Data

Other Research Data

Survey Data

## Genetic Data

GWAS

## Biobank Data

Samples
- DNA
- Serum
- Plasma

Consent
- Recontact
- Consent Status

## Informatics Tools

Calculated Controls (Charlson Index)

| Data Visualization | Data Queries |
| Annotation | Extract Data |

Natural Language Processing

### Validated Phenotypes

| Type II Diabetes | IBD |
| Coronary Artery Disease | Multiple Sclerosis |
| Congestive Heart Failure | Bipolar Disorder |
| Rheumatoid Arthritis | |

## Research

# Curating a Disease Algorithm

1. **Create a gold standard training set**.



2. **Create a comprehensive list of features** from patient's electronic data that describe the disease of interest



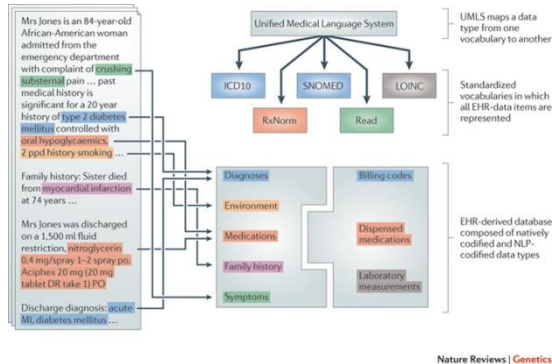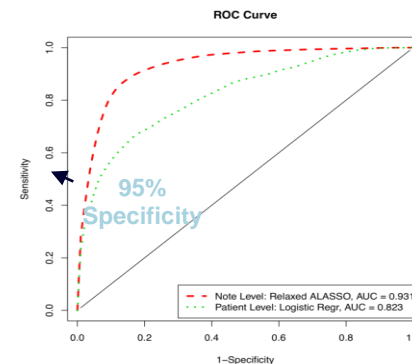Nature Reviews | Genetics

3. **Develop the classification algorithm**. Using the data analysis file and the training set from step 1, assess the frequency of each variable. Remove variables with low prevalence. Apply adaptive LASSO penalized logistic regression to identify highly predictive variables for the algorithm



4. **Apply the algorithm to all subjects** in the superset and assign each subject a probability of having the phenotype



95% Specificity

# Biobank Portal | Curated Diseases

| Validated Phenotype | Count* | Predictive Positive Value |
| --- | --- | --- |
| Bipolar Disease | 71 | 89% |
| Congestive Heart Failure | 387 | 90% |
| Coronary Artery Disease | 2,420 | 97% |
| Crohn's Disease | 453 | 90% |
| Multiple Sclerosis | 94 | 90% |
| Rheumatoid Arthritis | 550 | 90% |
| Type 2 Diabetes Mellitus | 1,887 | 97% |
| Ulcerative Colitis | 330 | 90% |

| Healthy Controls based on Charlson Index | Count** |
| --- | --- |
| 0 – 10-year survival probability is >98.3% | 2,206 |
| 1 – 10-year survival probability is >95.87% | 4,343 |
| 2 – 10-year survival probability is >90.15% | 6,545 |

* Based on 15,880 patients
** Based on 21,300 patients

# Automated Learning Algorithms enabled in RPDR such as PheNorm Algorithm

**Enabling phenotypic big data with PheNorm.**

Yu S[1,2], Ma Y[3], Gronsbell J[4], Cai T[5], Ananthakrishnan AN[6], Gainer VS[7], Churchill SE[8], Szolovits P[9], Murphy SN[7,10], Kohane IS[8], Liao KP[11], Cai T[4].

Machine Learned Query Terms
- Phenotypes New
  - Cardiology
  - Gastroenterology
  - Hematology
  - Metabolic diseases
  - Neurology
    - Cerebral aneurysm (IA)
    - Epilepsy (EPIL)
    - Insomnia (INSOM)
    - Intracranial hemorrhage (ICH)
    - Ischemic stroke (ISTR)
    - Migraine headache (MHA)
    - Multiple sclerosis (MS)
    - Obstructive sleep apnea (OSA)
    - Parkinson's disease (PD)
      - Current or Past History of PD
  - Oncology
  - Psychiatry
  - Pulmonology
  - Rheumatology
  - Urology

$$z = \log(1 + x) - \alpha * \log(1 + x_{note})$$

**Step 1: Normal Mixture Normalization**

Raw Feature

Normalized Feature

Accuracy Improvement

$$y \sim \tilde{Z}$$

**Step 2: Random Corruption Denoising**

# Machine Learned Phenotypes

- Abdominal hernia
- Acute bronchitis and bronchiolitis
- Acute pancreatitis
- Alcoholism
- Alzheimer's disease
- Aortic aneurysm
- Aplastic anemia
- Atrial fibrillation
- Atrioventricular block
- Autism spectrum disorders
- Basal cell carcinoma
- Bipolar Disease
- Bladder cancer
- Brain cancer
- Breast cancer
- Cerebral aneurysm
- Cholelithiasis
- Chronic pancreatitis
- Chronic sinusitis
- Coronary atherosclerosis
- Crohn's disease
- Deep vein thrombosis
- Depression
- Diverticulosis and diverticulitis
- Eating disorder
- Epilepsy
- Gastroesophageal reflux disease
- Gout
- Heart valve disorders

- Hyperlipidemia
- Hyperparathyroidism
- Hypertension
- Hypothyroidism
- Insomnia
- Intracranial hemorrhage
- Ischemic stroke
- Leukemia
- Lung cancer
- Melanoma
- Migraine headache
- Multiple sclerosis
- Myocardial infarction
- Neutropenia
- Non-Hodgkin lymphoma
- Obesity
- Obsessive compulsive disorder
- Obstructive sleep apnea
- Ovarian cancer
- Pancreatic cancer
- Parkinson's disease
- Peripheral vascular disease
- Pneumonia
- Polycystic ovaries
- Prostate cancer
- Pulmonary heart disease
- Renal cancer
- Renal failure
- Schizophrenia
- Substance addiction

- Suicidal ideation
- Suicide attempt or self-inflicted injury
- Thyroid cancer
- Tobacco use disorder
- Type 1 diabetes
- Type 2 diabetes
- Ulcerative colitis
- Urinary calculus
- Uterine cancer

# Phenotype Automation: Phenotype Quality Dashboard

# High Quality Phenotypes for Research Studies

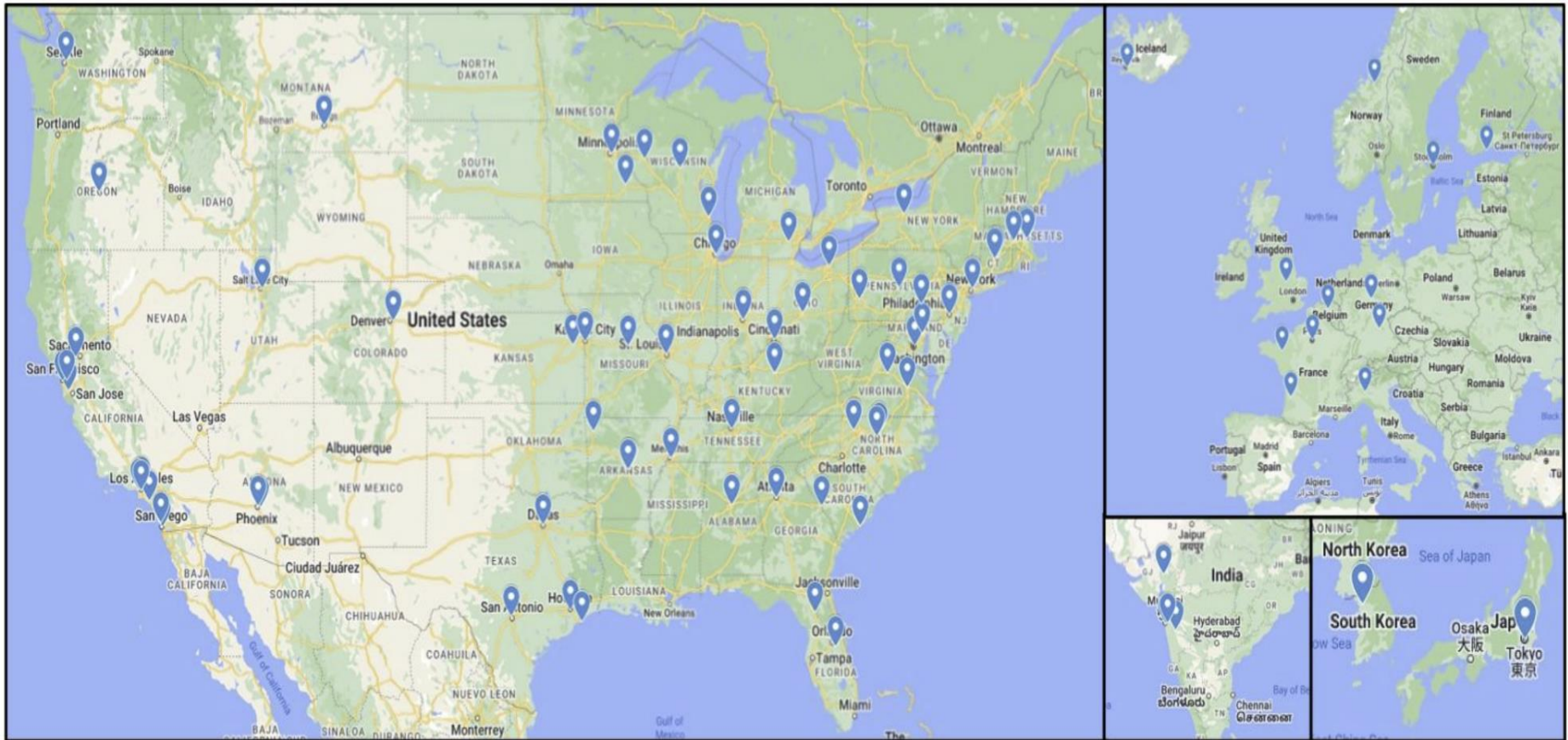# Combined with Generative AI can produce Digital Twin of Patient

# High Throughput Methods for supporting Translational Research

- Set of patients is selected from medical record data in a high throughput fashion

- Investigators explore phenotypes of these patients using Machine Learning tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Digital medicine is delivered into clinical care through Digital Twin

# I2b2 Implementations
>250 across the USA and Internationally, some illustrated below:

# Federated Queries



Mass General Brigham

Boston Children's Hospital

BIDMC

Boston Health Net (BMC and Community Health Centers)

Columbia U. Medical Center and New York Presbyterian Hospital

University of California, Davis

Washington University in St. Louis

Wake Forest Baptist Medical Center

Morehouse/Grady/RCMI

U Texas Health Science Center/Houston

# Drive Pragmatic Clinical Studies

# RECOVER Study Data Harmonization



https://recovercovid.org

# Data harmonized within i2b2 star schema

# Concepts in database available in harmonized ontology

# New i2b2 Query Tool to be released:

# I2B2 AI

## User asks a question



## Result rendered in web client



## AI returns response in i2b2 format

# AI-ENABLED QUERY BUILDER:
## *(I.E. INSTRUCTION-TUNED POC)*



i2b2

Export Query Log (XML)

i2b2 API

i2b2 XML Query

VAL

Validation

Web UI

Query Builder

XML INSTRUCT

Inference

Code Llama (LLaMa2)

LORA

4-Bit Quant

nmitchko/i2b2-querybuilder-34b-merged

i2b2 LLM POC Cell (@MGB)

Data Cleaning

Randomize

Tokenize

Data Prep

TRAIN

TEST

VAL

Query Validation

Code Llama built on Llama 2, trained on 500B tokens of *code*

Fine Tuning

Inference

Llama (LLaMa2)

LORA

4-Bit Quant

nmitchko/i2b2-querybuilder-34b-merged

Training Sandbox (@Mitchko Labs)

## High Throughput Methods for supporting Translational Research

- Set of patients is selected from medical record data in a high throughput fashion

- Investigators explore phenotypes of these patients using Machine Learning tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Digital medicine is delivered into clinical care through Digital Twin

# Congestive Heart Failure



- Affects 2% of the adult population

- Risk of death first year after diagnosis:  35%

- In patient hospital costs in 2011: $10.5B which is a small fraction of all heart failure related care

# Early Detection of Worsening or Improving Anemia

## Background and Methods

- Anemia is one of the strongest predictors of morbidity and mortality in CHF.

- Increasing or decreasing HGB is a further strong predictor, but there is no good way to determine whether a patient's HGB is on its way up or down (*Circulation. 2005;112:1121-1127*)
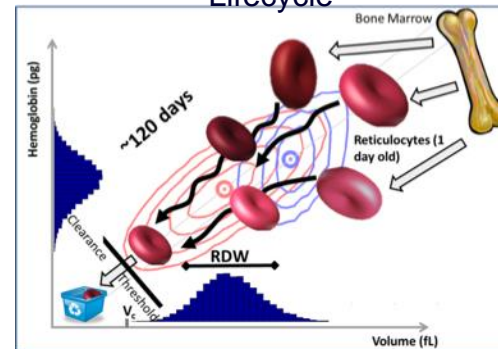
## Results and Conclusions

- A novel mathematical model of the RBC lifecycle enables estimation of patient-specific rates of RBC maturation and turnover from a routine CBC.
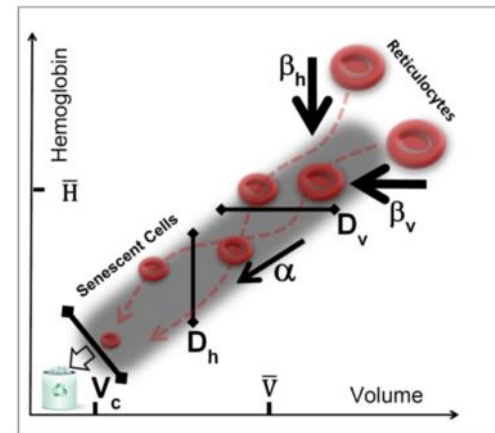
## Applications

1. CHF patients most likely to have decreasing HGB may benefit from altered treatment or longer hospitalization to avoid readmission.

2. CHF patients most likely to have increasing HGB may be responding well to treatment and benefit from earlier discharge or maintenance of current therapy.
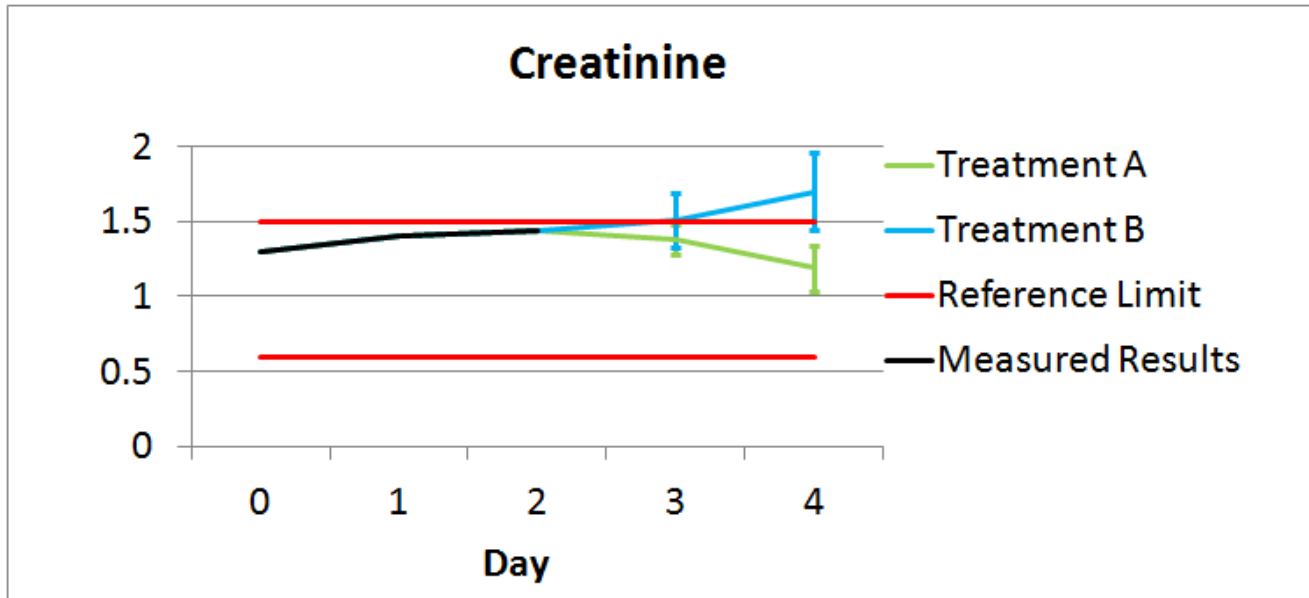
Dynamic Model of the RBC Lifecycle



Quantify Maturation and Clearance Rates

# Creatinine Prediction: Hypothetical Application



- Hypothetical analysis of creatinine times series where possible treatments are introduced into the model

- The model hypothetically provide a future trajectory conditioned on each treatment

# Bringing Big Data into Clinical Care with Open App Development



**Clinician**

**SMART App embedded in Epic**

Epic Data Repository

DATA
GeneInsight, mHealth, ePath, Medical Images, 25 years of Legacy electronic data, and Other External Systems

Analytic Calculation Engine

Analytics have direct access to repository

Core Integration Database

FHIR interface for real time updates

**Laboratory Personnel**

**SMART App in Lab**

Non-EHR Users View Standalone App

# Transforming Care in the Digital Age



**Digital and IoT devices continuously output Patient Data**

## PATIENT

**Digital Twin of patient enables continuous assessment of patient with Real Time Algorithms**

**Navigator Model dramatically increases Frequency and Convenience for Patient Communication**
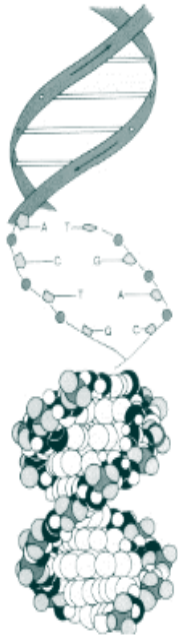
**System drives Pragmatic Clinical Trials Leading to Continuous Process Improvement**

# MGB Data Enclave Overview

# I2b2 tranSMART Software

*i2b2 Homepage (https://www.i2b2.org)*

*i2b2 Software (https://www.i2b2.org/software)*

*i2b2 Community Site (https://community.i2b2.org)*

*https://i2b2transmart.org/2023-i2b2-symposium/2023-symposium-recordings-slides/*