

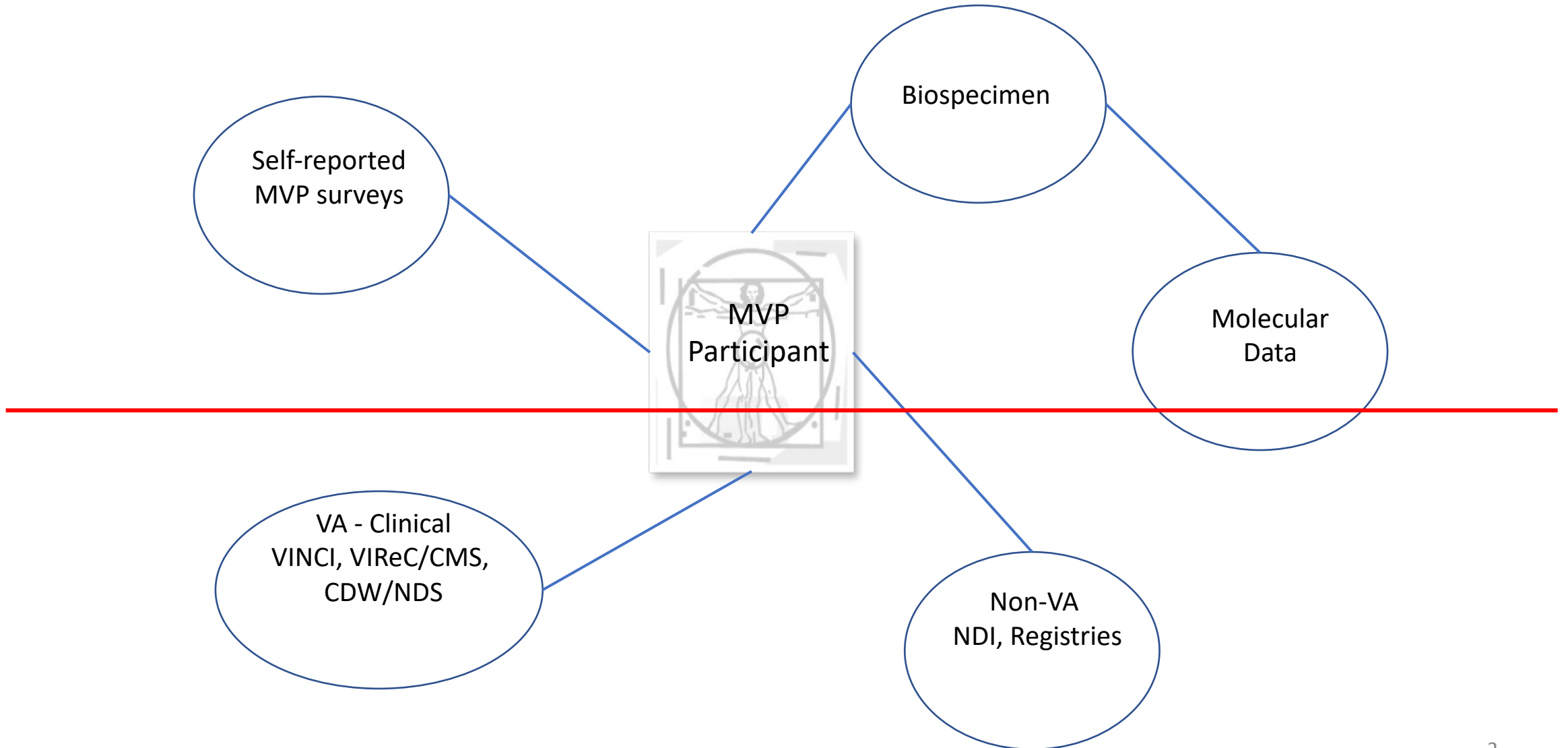
Precision Medicine: Understanding How Genomics and Big Data Are Enabling the Right Care for the Right Person at the Right Time

J. Michael Gaziano, MD, MPH

**Chief, Division of Aging, Brigham and Women's Hospital
Director, MAVERIC; PI, MVP, VA Boston Healthcare System
Professor of Medicine, Harvard Medical School**



Organizing the MVP Data Universe



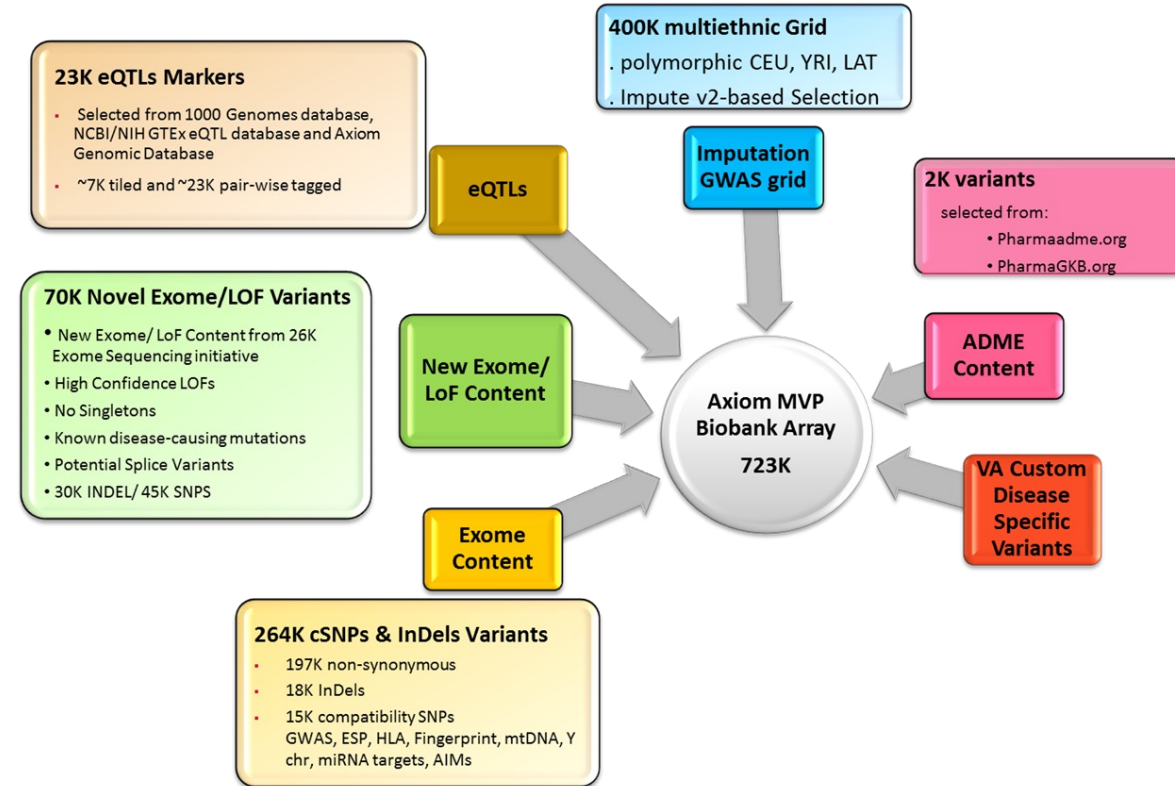
Specimen Stats and Omics Strategy

- **VA Central Biorepository**

- Over 1,000 samples/day processed/DNA isolated
- Over 8 million aliquots at -80°C (Boston; Albq)
- Send out over 200,000 per year

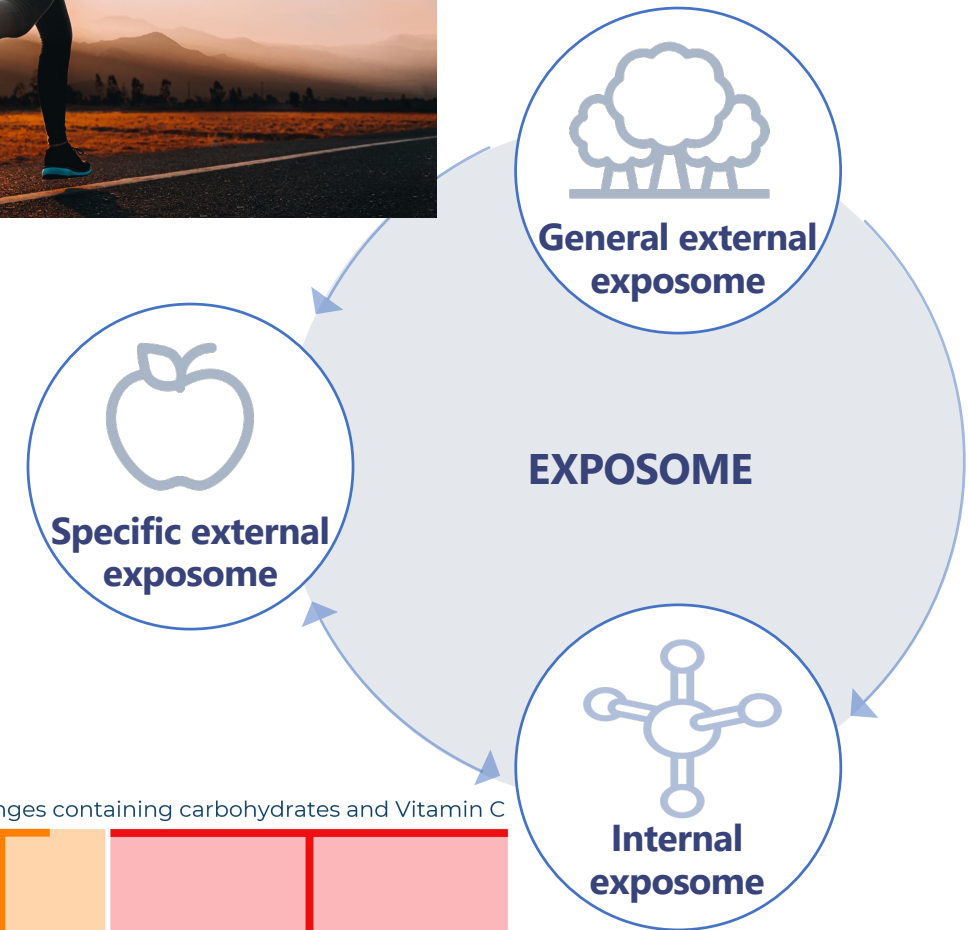
- **Aliquots shipped 2011 – 2023**

- DNA
 - Genotyping: 943,274
 - Ethnic Focus Array: 256,657
 - Whole Genome Sequencing (30x): 186,085
 - Methylation: 94,020
- Plasma
 - Metabolomic: 60,160
 - Lipidomic: 40,000
 - Proteomic: 1,709
- Future: other collections, plasma-based RNA fragments, echo of microbiome

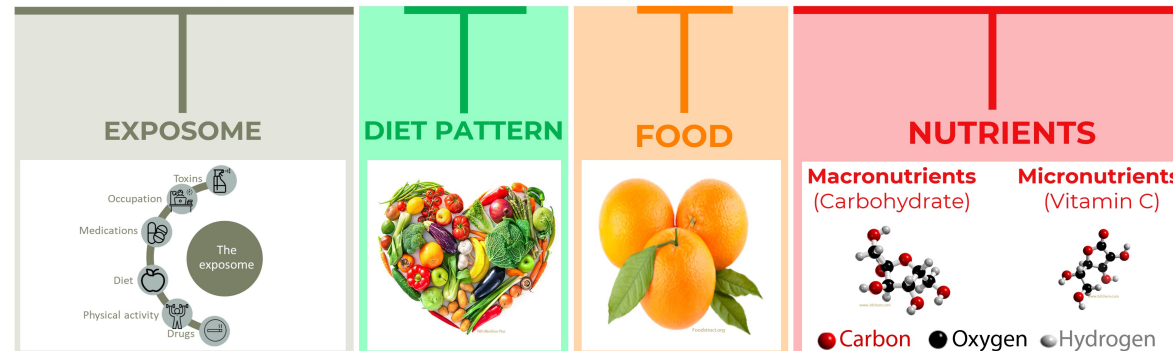


Exposome

- Lifestyle
 - Smoking
 - Physical activity
 - Diet/Alcohol
- Environmental factors
 - Geospatial data
 - Military exposure
- Social deprivation
- Others



A Veteran who enjoys running, is a vegetarian who eats oranges containing carbohydrates and Vitamin C



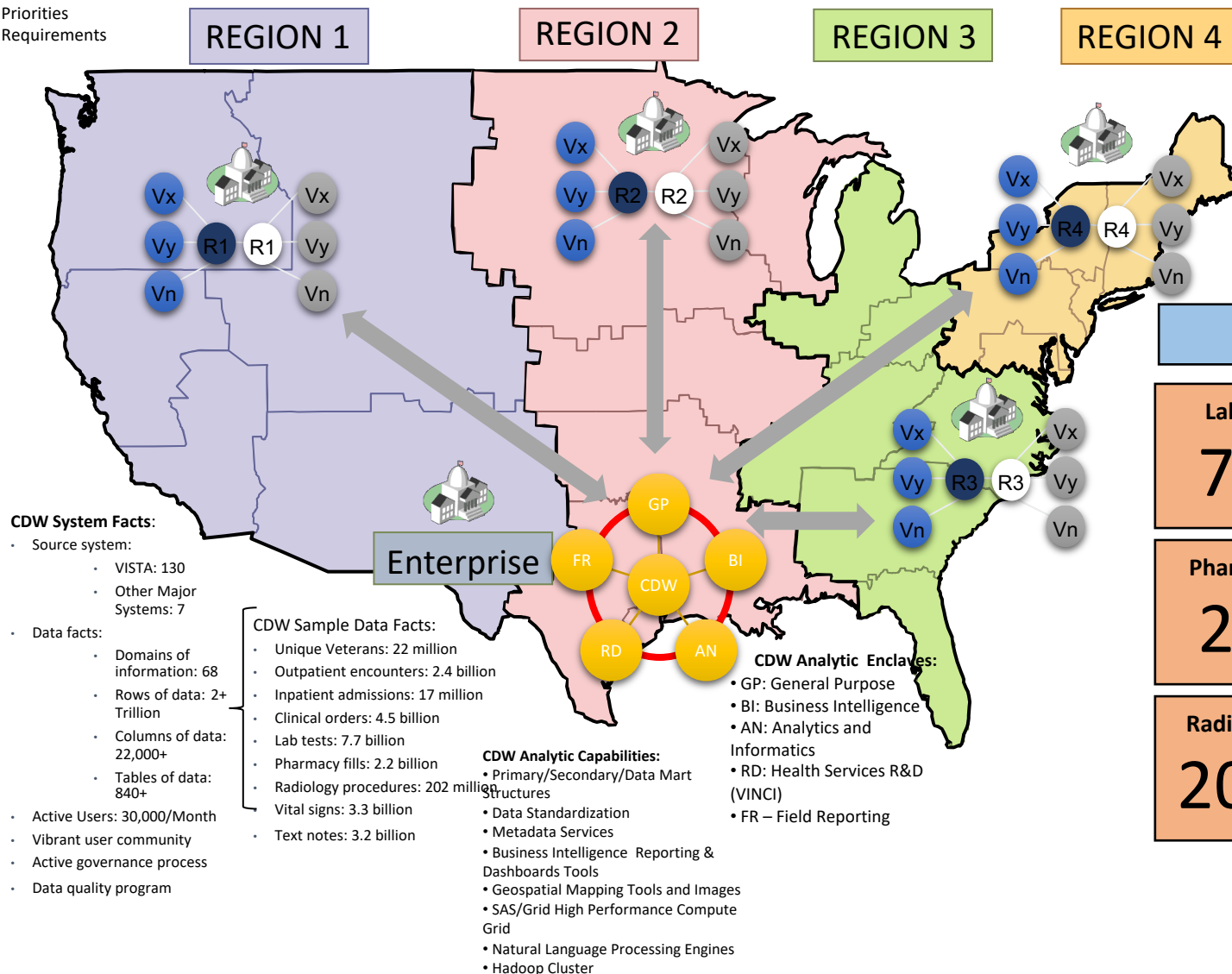


VA Analytic Ecosystem (2015)

Common Data ♦ Common Infrastructure ♦ Common Tools ♦ Common Security

Governance Board

- Strategy
- Policy
- Priorities
- Requirements



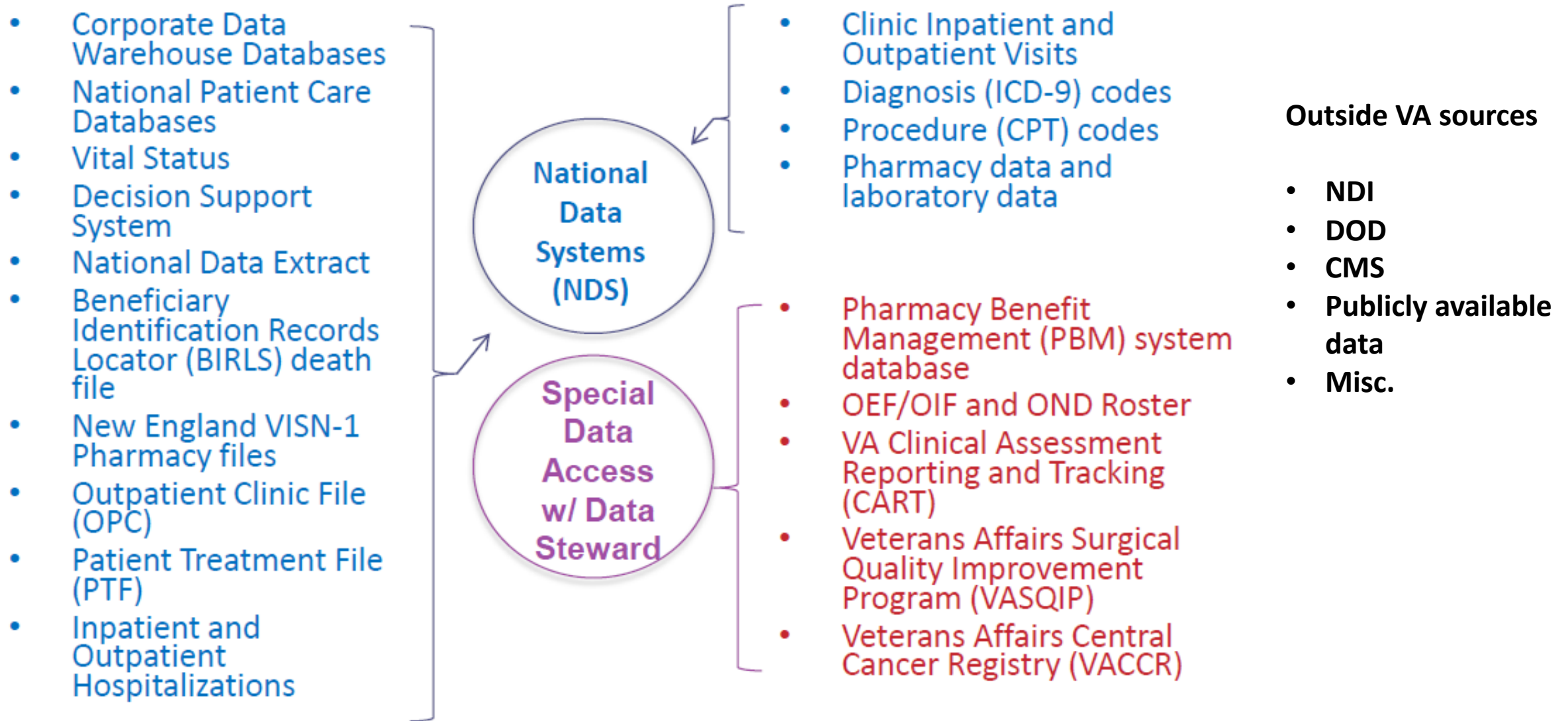
The EHR data available in the CDW and other data sources is among the best in the world.

Patients: 22 M

Lab Results 7.7B	Clinical Orders 4.5B	Immunizations 71M	Appointments 1.4B
Pharmacy Fills 2.2B	Clinical Notes 3.2B	Health Factors 2.2B	Encounters 2.4B
Radiology Proc 202M	Vital Signs 3.3B	Consults 315M	Admissions 17M
Surgeries 14M		Oncology 1.3M	

Domains: 15/68

MVP VA and Other Data Sources



General Phenotyping Goals

More and more data are becoming available for research:
Is it a blessing or a curse?

- Opportunities and challenges
- Are there appropriate tools and resources to analyze, manage and handle these data?
- Are we optimally synthesizing all the information?
- Do we have all the information and annotation?

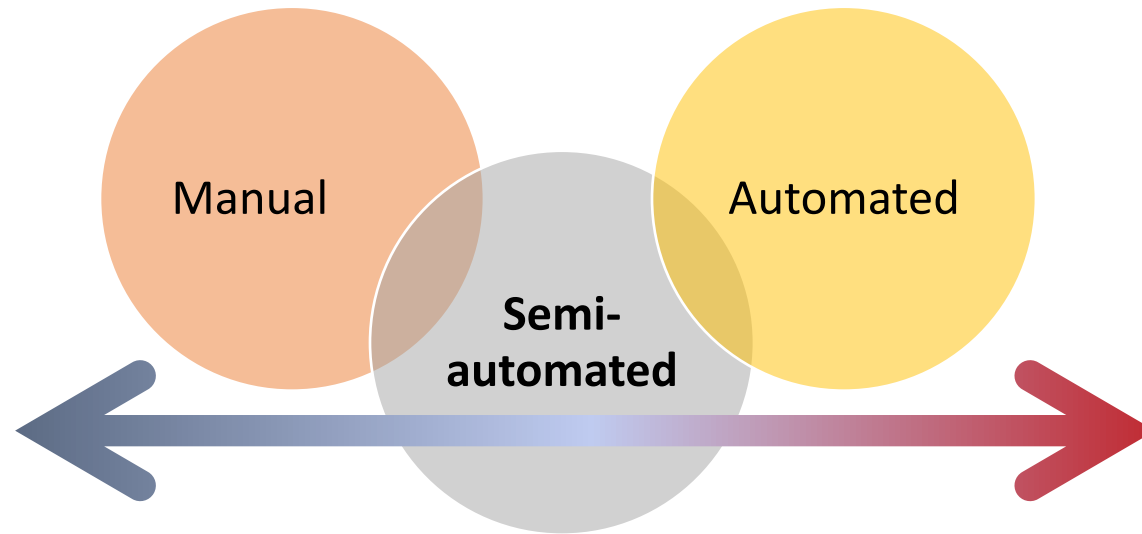


Sometimes, data warehouses resemble landfills more than libraries.

-Phenotypes are the foundation of clinical research

-Major challenge is in accurately and efficiently assigning phenotypes to subjects

Our Vision for Advanced Phenotyping in MVP: A New Approach



Semi-automated phenotyping combines features
of manual and automated phenotype development

Manual Curation: Laboratory Adjudication Effort

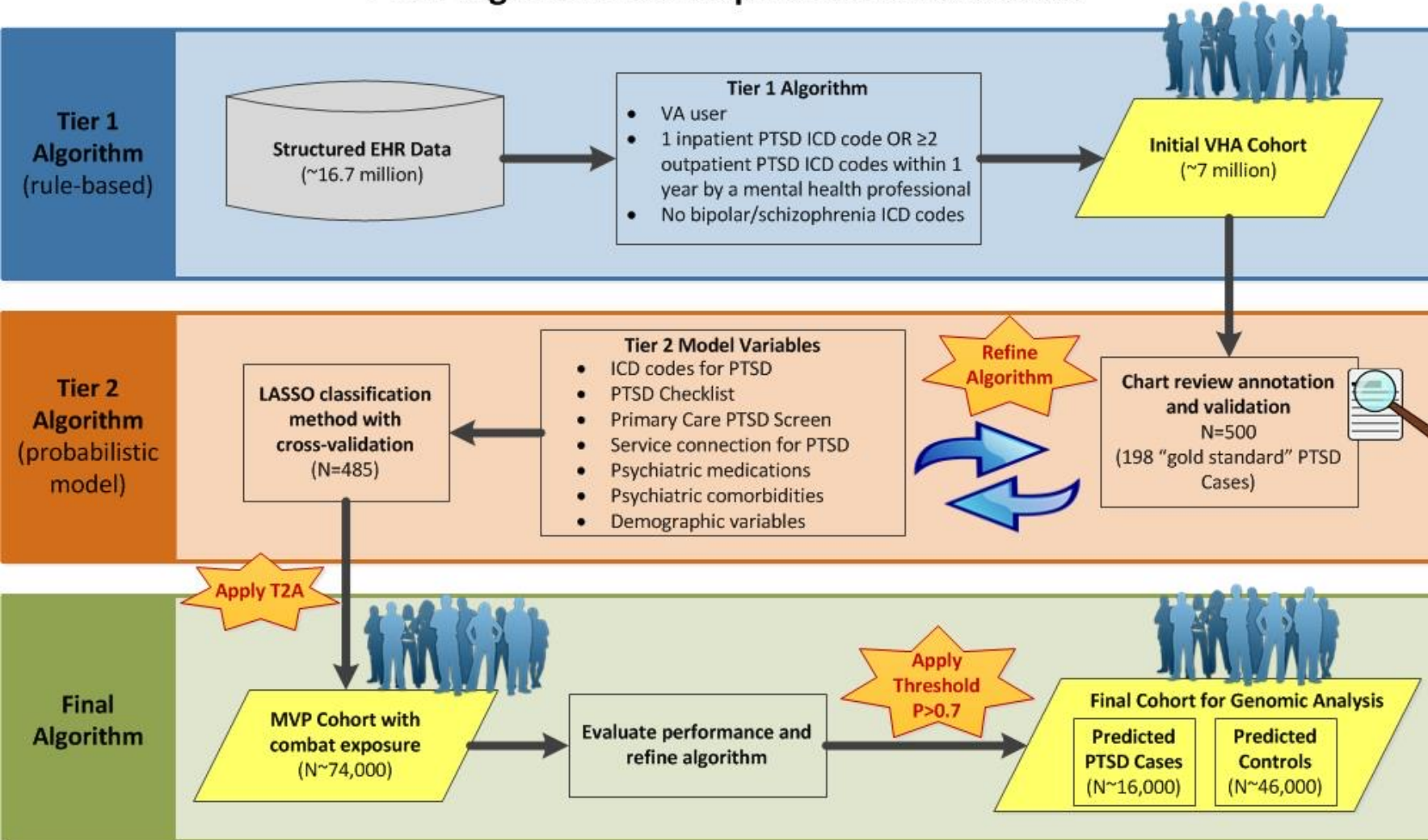
Laboratory test name	Number of tests adjudicated	Number of tests accepted
Hemoglobin A1C	527	365
Serum albumin	4141	644
Blood Glucose	4578	905
HDLC	770	377
Hemoglobin	2638	331
LDLC	1230	602
Serum Potassium	2198	720
Serum Creatinine	5212	705
Serum Sodium	2608	757
Total Cholesterol	2137	405
Triglycerides	1528	390

Serum Albumin Adjudication

Accept	LabChem TestSID	LabChem TestName	Specimen	VISN	Sta3n	Units	n	min	p1	p5	p10	p25	p50	p75	p90	p99	max
Yes	800000948	ALBUMIN(SEATTLE)	Serum	20	648	G/DL	8985	-0.22	3.1	3.7	3.9	4.2	4.4	4.6	4.8	5.2	6
No	800001031	albumin(ep), csf	Cerebral spinal fluid	20	648	%	22	51	51	54	55	57	61	66	69	71	71
No	800001092	MICROALBUMIN	Urine	20	648	MG/DL	70167	0	0.3	0.43	0.7	1.28	2.8	8.28	30.4	228.8	21321
Yes	800001119	ALBUMIN	Plasma	20	648	g/dL	712338	0.1	1.9	2.6	3.1	3.8	4.2	4.4	4.6	5	67
Yes	800001119	ALBUMIN	Serum	20	648	g/dL	21999	0.2	2.1	2.7	3.2	3.9	4.3	4.5	4.7	5.1	7.6

Manual Curation: PTSD Phenotype – CSP 575B/MVP

PTSD Algorithm Development and Validation



Purpose: To develop and validate EMR-based algorithm for identifying PTSD in a sample of Veterans using a probabilistic modeling approach

Validation of an Electronic Medical Record-Based Algorithm for Identifying Posttraumatic Stress Disorder Cases in a VA Million Veteran Program Sample Using a Multi-Tiered Phenotyping Approach. Kelly M. Harrington, Rachel Quaden, Jacqueline Honerlaw, Murray Stein, Joel Gelernter, Shadha Cissell, Robert Pietrzak, Krishnan Radhakrishnan, John Michael Gaziano, John Concato, David R. Gagnon †, and Kelly Cho † on behalf of the VA Million Veteran Program

A phenotyping algorithm to identify acute ischemic stroke accurately from a national biobank: the Million Veteran Program

This article was published in the following Dove Press journal:
Clinical Epidemiology

Tasnim F Imran,^{1,3,*} Daniel Posner,^{1,4,*} Jacqueline Honerlaw,¹ Jason L Vassy,^{1,2} Rebecca J Song,¹ Yuk-Lam Ho,¹ Steven J Kittner,⁵ Katherine P Liao,^{1,2} Tianxi Cai,^{1,6} Christopher J O'Donnell,^{1,2} Luc Djousse,^{1,2} David R Gagnon,^{1,4} J Michael Gaziano,^{1,2} Peter WF Wilson,^{7,8} Kelly Cho^{1,2}
On behalf of the VA Million Veteran Program

¹Massachusetts Veterans Epidemiology Research and Information Center (MAVERIC), VA Cooperative Studies Program, VA Boston Healthcare System, Boston, MA, USA; ²Department of Medicine, Division of Aging, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA; ³Department of Medicine, Cardiology Section, Boston

Background: Large databases provide an efficient way to analyze patient data. A challenge with these databases is the inconsistency of ICD codes and a potential for inaccurate ascertainment of cases. The purpose of this study was to develop and validate a reliable protocol to identify cases of acute ischemic stroke (AIS) from a large national database.

Methods: Using the national Veterans Affairs electronic health-record system, Center for Medicare and Medicaid Services, and National Death Index data, we developed an algorithm to identify cases of AIS. Using a combination of inpatient and outpatient ICD9 codes, we selected cases of AIS and controls from 1992 to 2014. Diagnoses determined after medical-chart review were considered the gold standard. We used a machine-learning algorithm and a neural network approach to identify AIS from ICD9 codes and electronic health-record information and compared it with a previous rule-based stroke-classification algorithm.

Results: We reviewed administrative hospital data, ICD9 codes, and medical records of 268 patients in detail. Compared with the gold standard, this AIS algorithm had a sensitivity of 91%, specificity of 95%, and positive predictive value of 88%. A total of 80,508 highly likely cases of AIS were identified using the algorithm in the Veterans Affairs national cardiovascular disease-risk cohort (n=2,114,458).

Conclusion: Our algorithm had high specificity for identifying AIS in a nationwide electronic health-record system. This approach may be utilized in other electronic health databases to accurately identify patients with AIS.

Imran TF, Posner D, Honerlaw J, Vassy JL, Song RJ, Ho YL, Kittner SJ, Liao KP, Cai T, O'Donnell CJ, Djousse L, Gagnon DR, Gaziano JM, Wilson PW, Cho K.

Clin Epidemiol. 2018 Oct 16;10:1509-1521. doi: 10.2147/CLEP.S160764.

Neural Network

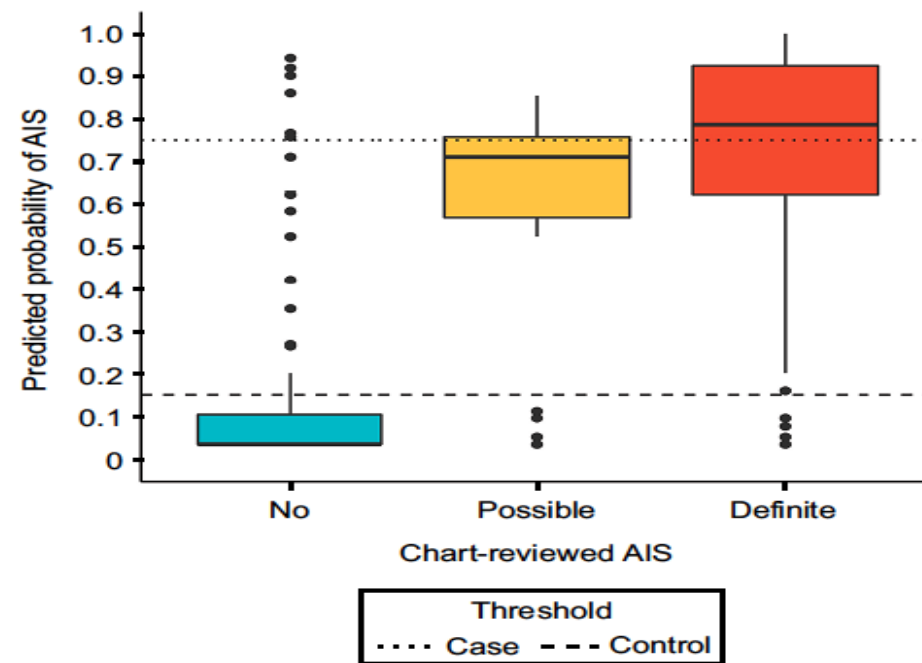
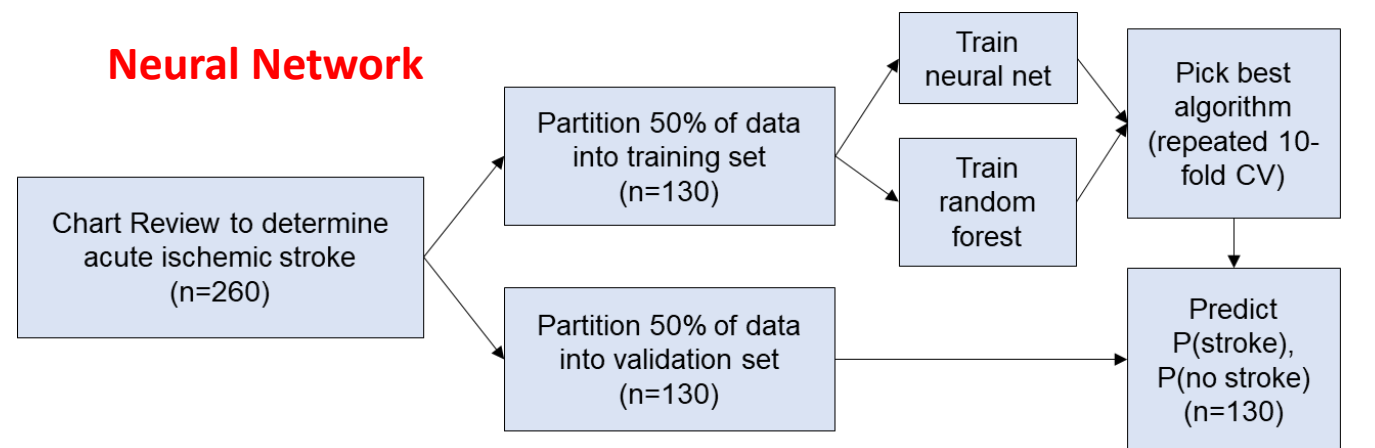


Figure 3 Predicted probabilities of stroke based on charts reviewed.

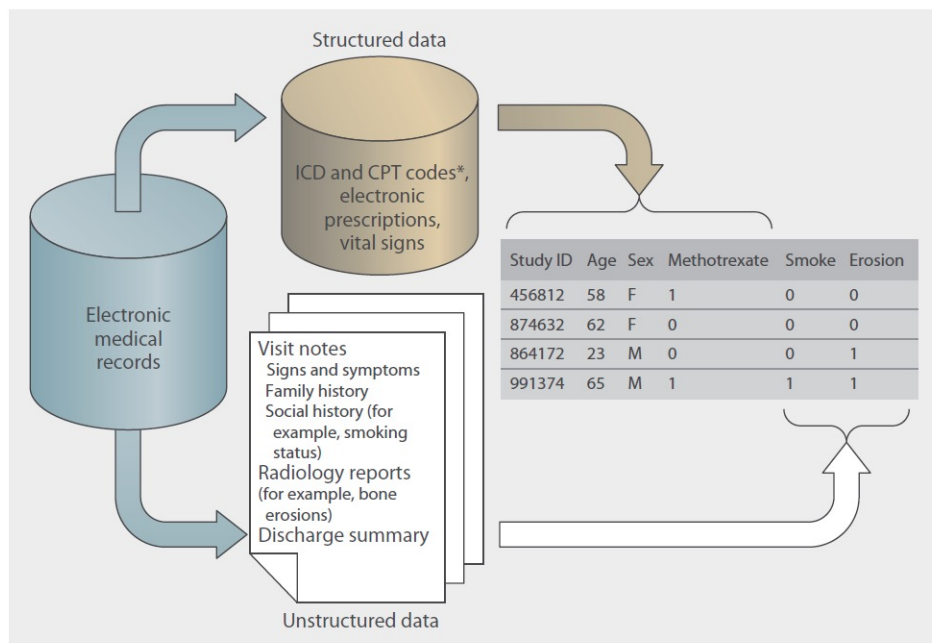
Notes: Thresholds optimized for largest n (excluding $P_{\text{control}} < P < P_{\text{case}}$) with Cohen's $\kappa > 0.9$ between algorithm labels and review labels.

Abbreviation: AIS, acute ischemic stroke.

MAVERIC

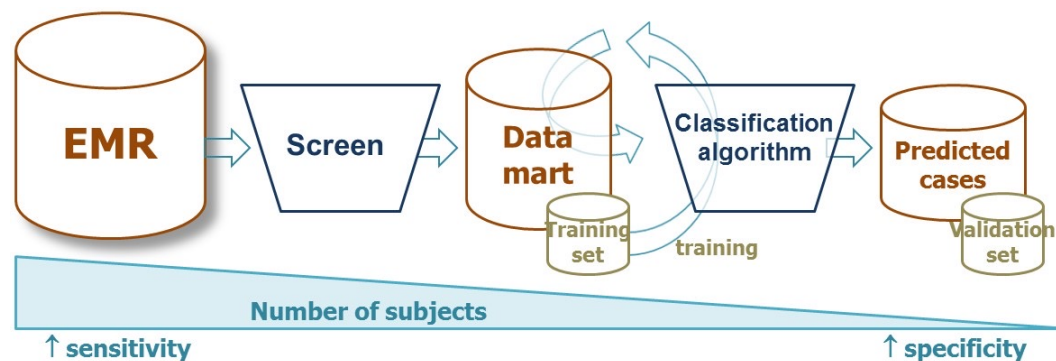
Massachusetts Veterans Epidemiology Research and Information Center

Advanced Phenomics



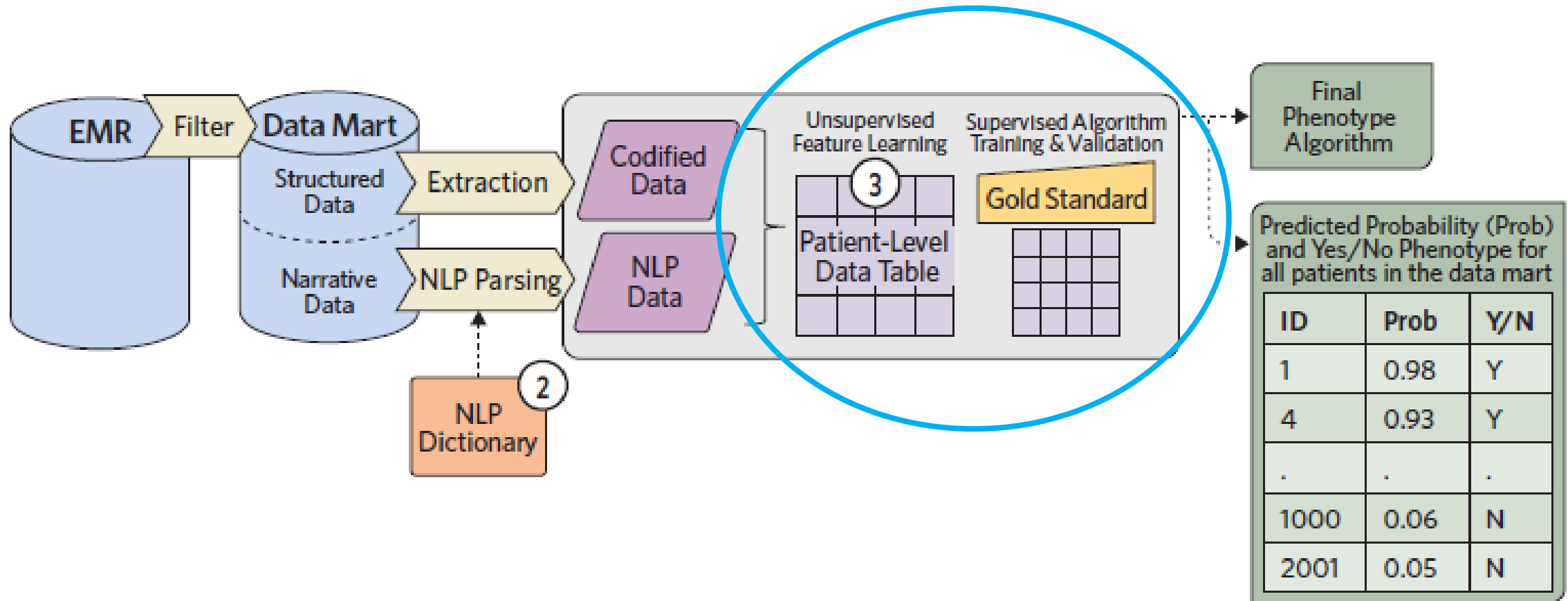
This block contains two screenshots of web resources and a central hexagonal overlay representing NLP tasks. The left screenshot is from Wikipedia, showing the article for 'Rheumatoid arthritis'. The right screenshot is from Medscape, showing the 'Rheumatoid Arthritis' reference page. Below these is a screenshot of the UMLS (Unified Medical Language System) website. Overlaid on the right side of these screenshots is a cluster of six blue hexagons, each containing a text processing task:

- Term Detection
- Concept Mapping
- Drug Grouping
- Junk Filtering
- Frequency Control
- RankCor Control

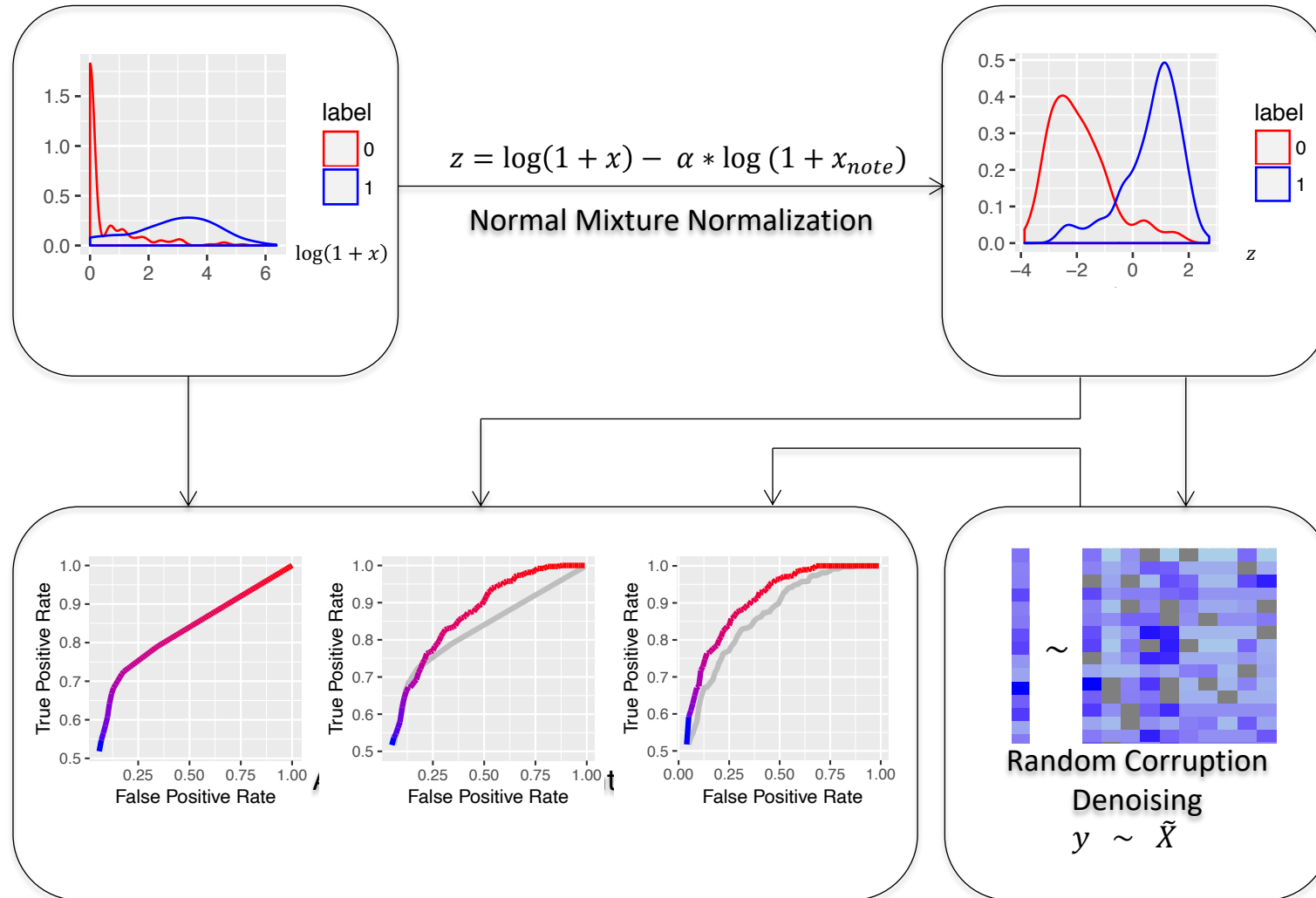


Semi-supervised: Machine learning, NLP, and EHR

Pipeline for phenotyping (PheCAP)



Unsupervised: PheNorm workflow

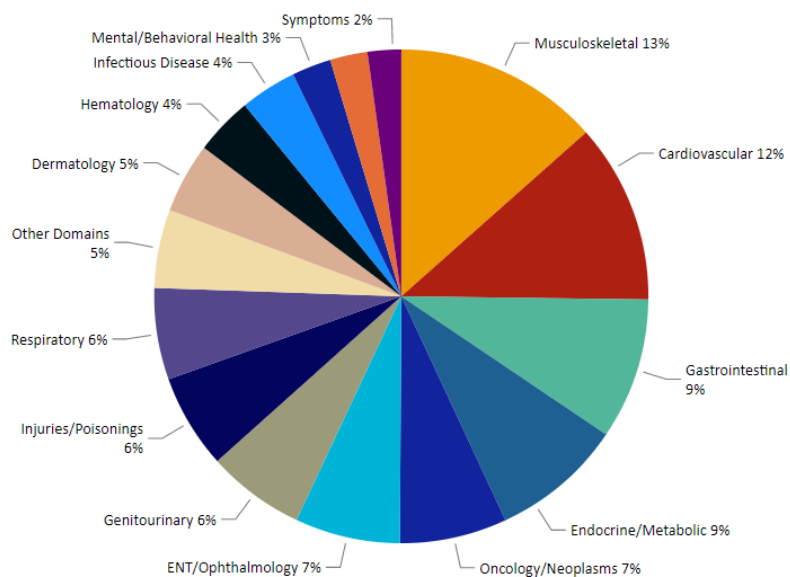




Total Phenotypes	6,000+
Total user base	7,400+
CIPHER Online site launch	June 2023



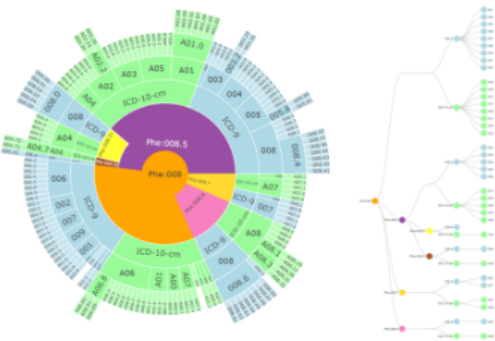
<https://phenomics.va.ornl.gov/>



Note: Phenotypes may fall into more than one data domain

Data Visualization Tools

This page contains tools to allow users to visualize expand the details section under each instrument



CIPHER Online

<https://phenomics.va.ornl.gov/>

The screenshot shows the 'General Phenotype' entry form. It includes a header with the VA logo and 'CIPHER' text. A navigation bar has links for Home, Getting Started, Explore, About, and Contact Us. The main content area is titled 'General Phenotype' and contains a step indicator for '1 Basic Information and Contact'. Below this, there are input fields for 'Phenotype Name*', 'Abbreviations and Keywords*', and 'Author(s)*', each with a search bar and instructions. A 'Please complete all fields below. You will have the opportunity to review all information at the final step.' message is displayed at the top of the form.

Interactive Phenotype
Entry Form Wizard

The screenshot shows the 'Explore' page in CIPHER. It features a sidebar with filters for Data Classification, Related Disease Domain, Data Sources Used, Algorithm components, Role of phenotype in analysis, Date algorithm created, Author, Method used, Publication, Algorithm code, Validated, and Attachment. The main content area displays a list of phenotypes, including 'ASCVD (MAP)', 'Abdominal aortic aneurysm (MAP)', 'Abdominal hernia (MAP)', and 'Abdominal pain (MAP)', each with details like Author (MVP Core) and Algorithm Created date (10/31/2019).

Enhanced
Phenotype
Knowledgebase

The screenshot shows the 'Data Visualization Tools' page. It includes a header with the VA logo and 'CIPHER' text. A navigation bar has links for Home, Getting Started, Explore, About, and Contact Us. The main content area is titled 'Data Visualization Tools' and contains a description of the tools. Below the description, there are two visualizations: a circular diagram and a hierarchical tree diagram. To the right of the diagrams, there is a section titled 'Details' with sub-sections for 'OVERVIEW', 'USES', and 'AUTHOR & CITATION'.

Integrated Data
Visualization Tools

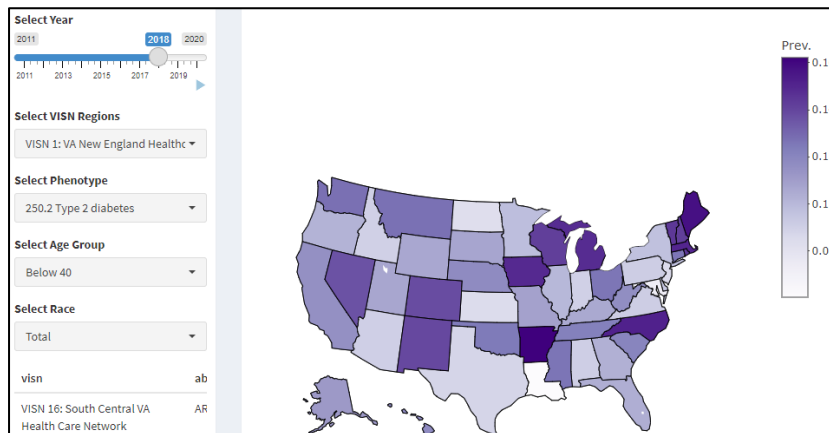
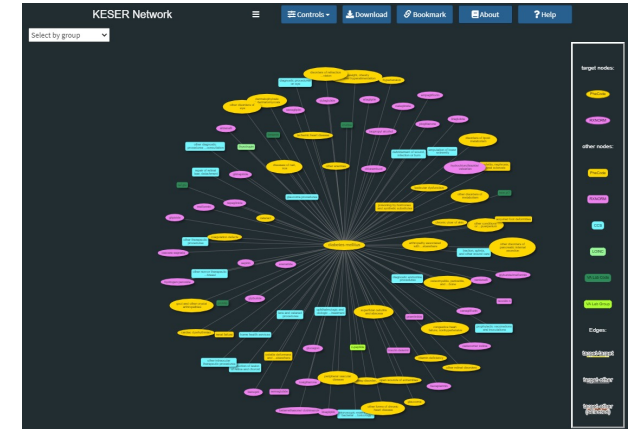
Data Visualization Tools

Data visualization tools integrated into phenotype definition knowledgebase

- KESER Network
- ICD Hierarchy Tool
- Geography of Phenotypes (GeoPheno)
- MVP GwPhewas

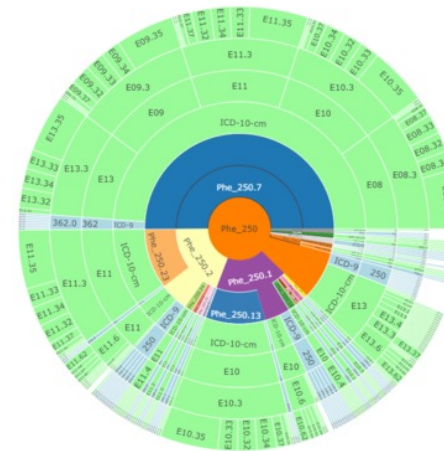
KESER Network*

Allows users to infer relatedness among diseases, treatment, procedures and laboratory measurements by creating a visual, interactive knowledge map



GeoPheno

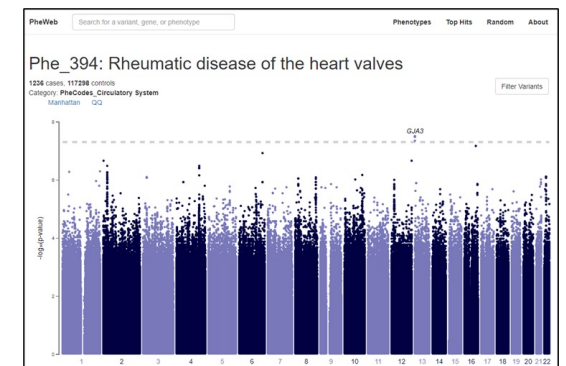
Display phenotype prevalence trends over time and location



ICD Hierarchy Tool*

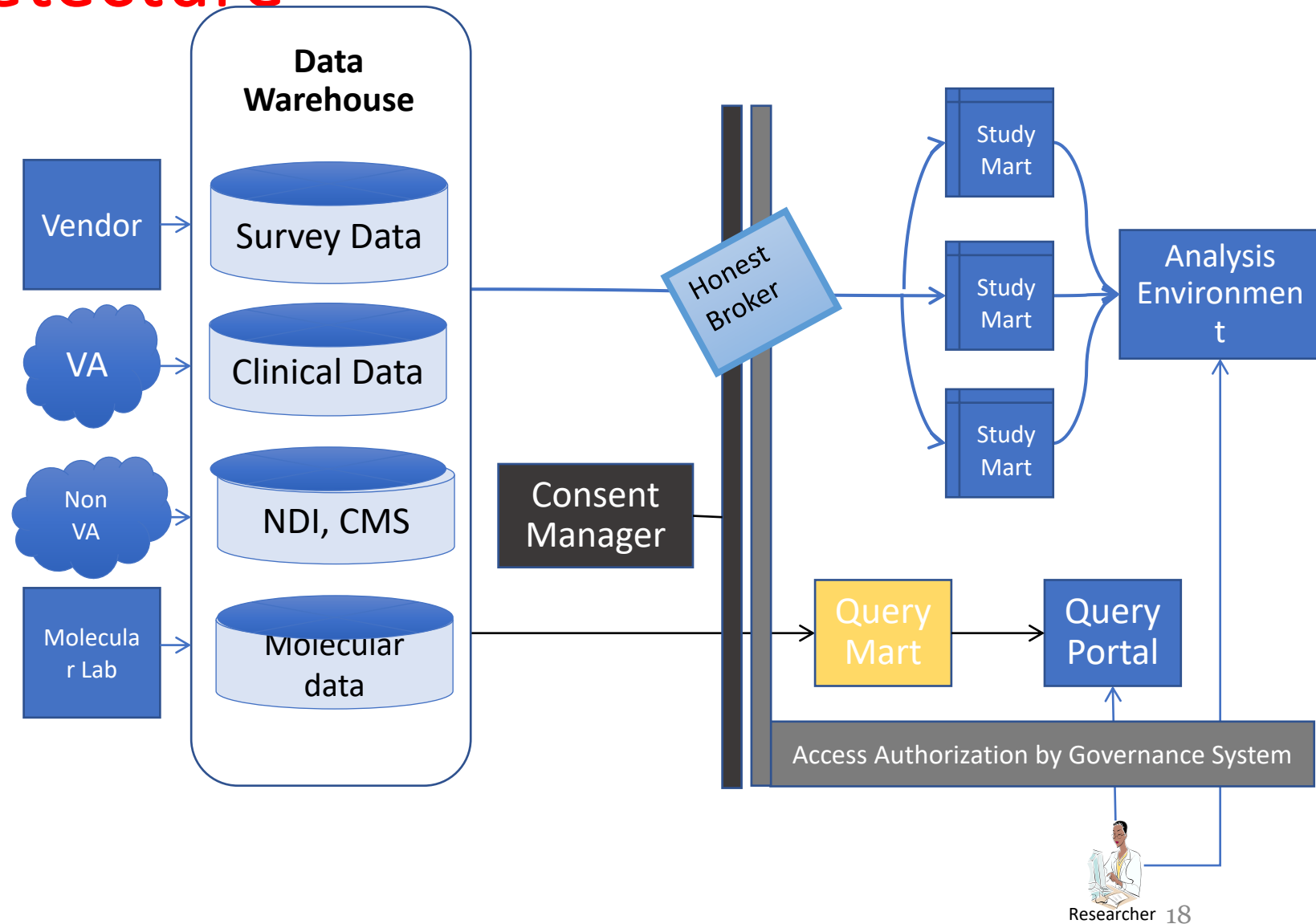
Enables searching across Phecode, ICD-9 and ICD-10 mappings

GwPhewas*
Display phenotype details and results summaries

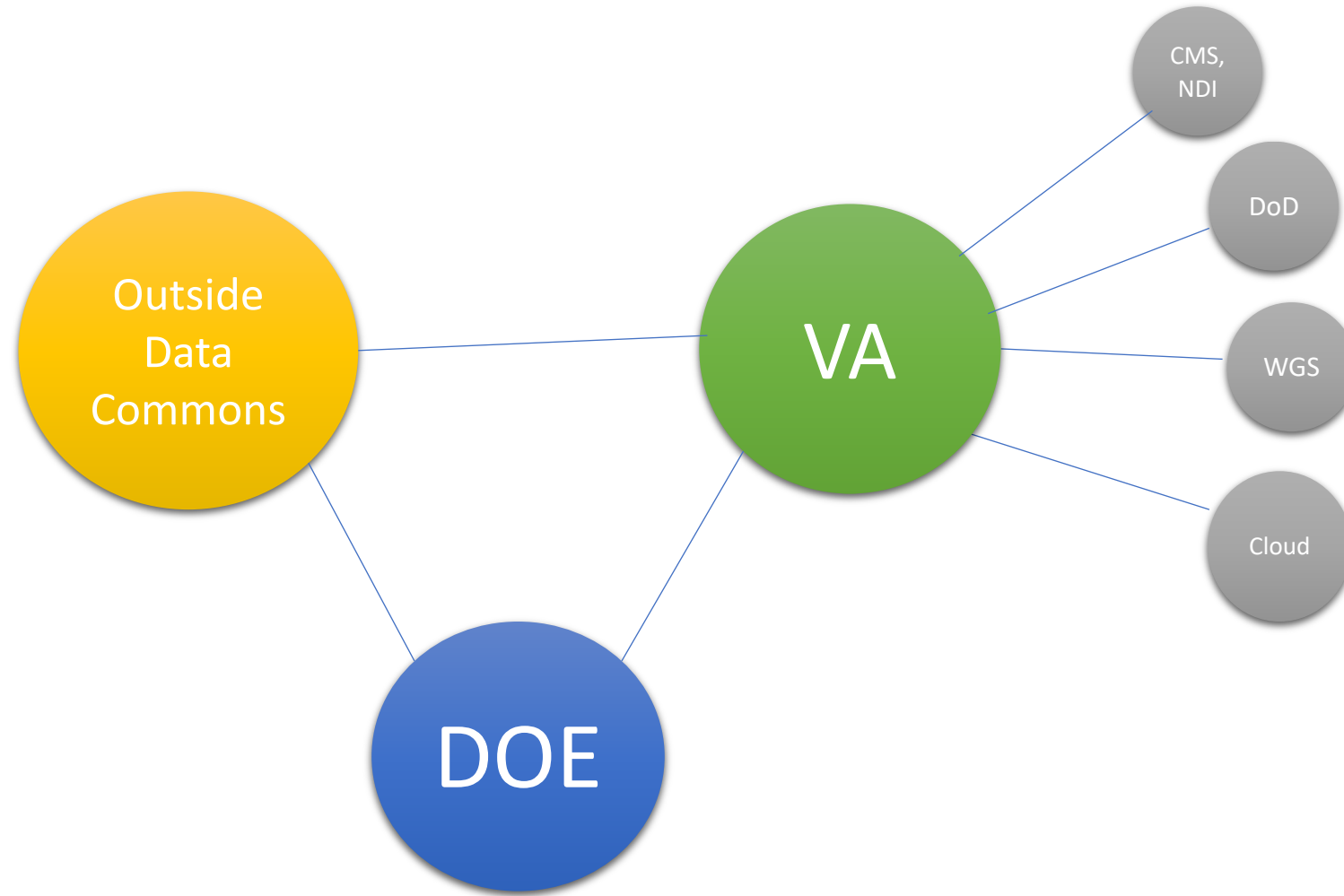


MVP Enterprise Architecture

- Requirements
 - Ingest phenomic and omic data from and/or have organic links to data sources.
 - Securely store data
 - Display data available
 - Provision data to many user types.
 - Track/monitor activity
- User types
 - Core users
 - Researchers
 - Industry partners
 - Others



MVP Computing Environment

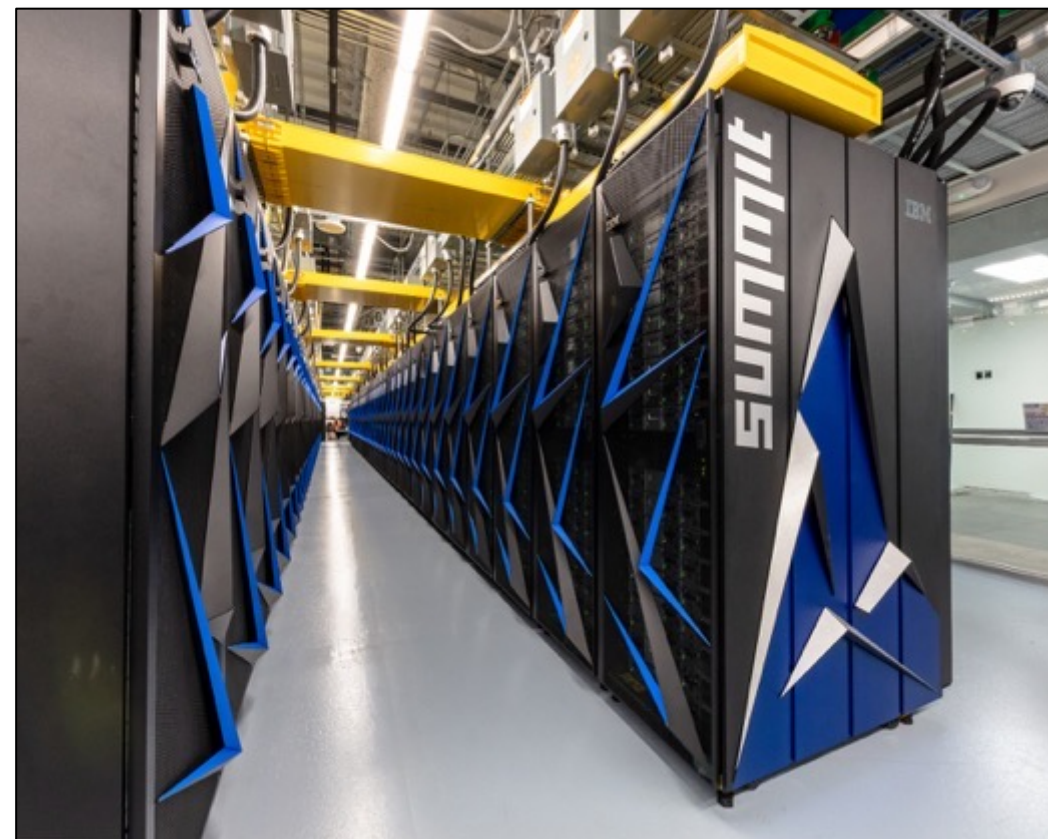




ORNL's computing vision enables a strong leadership position now and into the future

Sustain leadership and scientific impact in computing and computational sciences

- Provide world's most powerful open resources for scalable computing and simulation, data and analytics at any scale, and scalable cyber-secure infrastructure for science
- Follow a well-defined path for maintaining world leadership in these critical areas
 - Attract the brightest talent and partnerships
 - Deliver leading-edge science relevant to missions of DOE and key federal and state agencies
 - Invest in cross-cutting partnerships with industry
 - Provide unique opportunity for innovation based on multiagency collaboration
 - Invest in education and training



Using the MVP Data

MILLION
VETERAN
PROGRAM



Current MVP Studies

- Non-Genetic studies
- GWAS
- PRS
- Meta-analyses
- GWAS X PheWAS
- MR




MVP nutrition and exposome data behaves as expected

Dietary yogurt is distinct from other dairy foods in its association with circulating lipid profile: Findings from the Million Veteran Program

Kerry L. Ivey^{a b c}  , Xuan-Mai T. Ng,
Rebecca Song^{a m}, Geraint B. Rogers^{b f, 1},
Peter WF. Wilson^{g h}, Kelly Cho^{a d e}, J. Mi.
Walter C. Willett^{c j k}, Luc Djoussé^{a d e}

Article

The Structure of Relationships between the Human Exposome and Cardiometabolic Health: The Million Veteran Program

Kerry L. Ivey^{1,2,3,*} , Xuan-Mai T. Nguyen^{1,4,5},
Rebecca Song^{1,7}, Yuk-Lam Ho¹ , Ruifeng Li³ ,
John Michael Gaziano^{1,4,5,10}, Frank B. Hu^{3,11,12}, V



The American Journal of
CLINICAL NUTRITION


journal homepage: <https://ajcn.nutrition.org/>

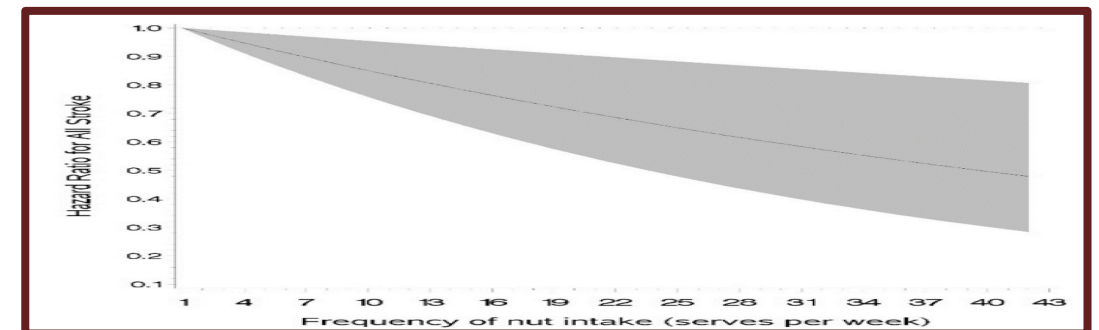
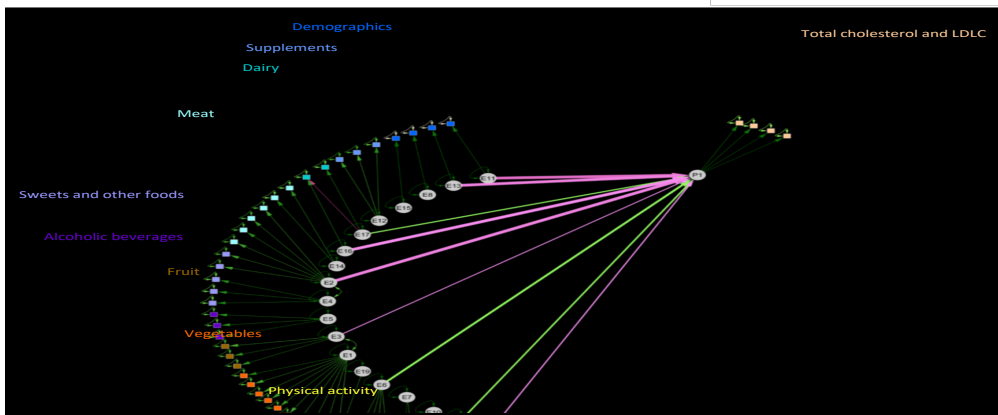
Original Research Article

Association of dietary fatty acids with the risk of atherosclerotic cardiovascular disease in a prospective cohort of United States veterans

Kerry L. Ivey^{1,2,3,*}, Xuan-Mai T. Nguyen¹, Ruifeng Li⁴, Jeremy Furtado⁴, Kelly Cho^{1,2,3},
John Michael Gaziano^{1,2,3,5}, Frank B. Hu^{4,6,7}, Walter C. Willett^{4,6,7}, Peter WF. Wilson^{8,9},
Luc Djoussé^{1,2,3,4}

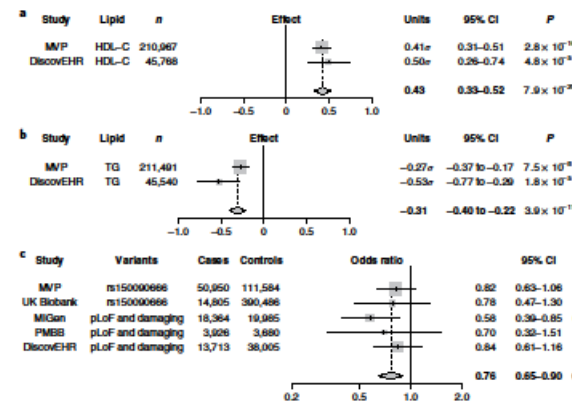
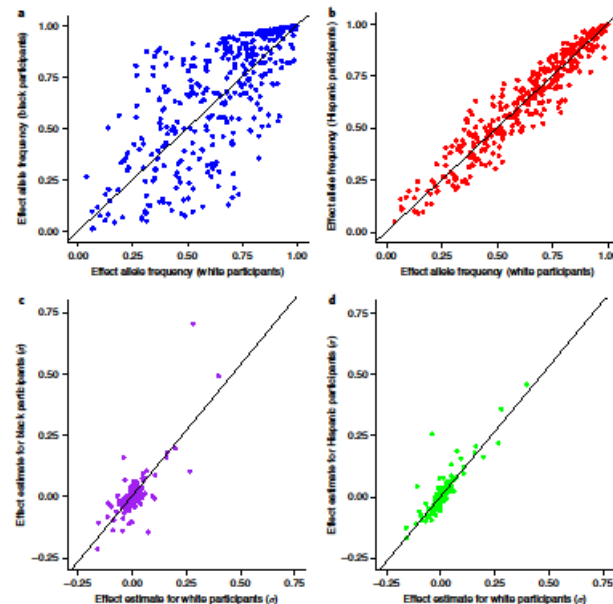
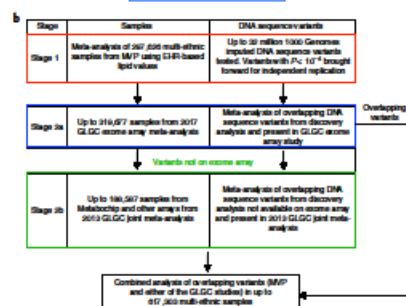
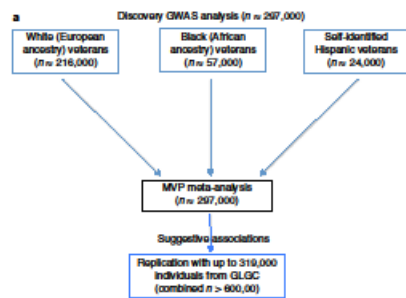
Association of Nut Consumption with Risk of Stroke and Cardiovascular Disease: The Million Veteran Program

Kerry L. Ivey^{1,2,3,*} , Xuan-Mai T. Nguyen^{1,4,5}, Rachel M. Quaden¹, Yuk-Lam Ho¹, Kelly Cho^{1,4,5},
J. Michael Gaziano^{1,4,5,6} and Luc Djoussé^{1,4,5,†}
















Genetics of blood lipids among ~300,000 multi-ethnic participants of the Million Veteran Program

Klarin, D, Damrauer SM, Cho K, et al, Global Lipids Genetics Consortium, Myocardial Infarction Genetics (MIGen) Consortium³¹, The Geisinger-Regeneron DiscovEHR Collaboration, The VA Million Veteran Program



Disease	Cases	Controls	Odds ratio	95% CI	P
Cardiovascular					
Coronary disease	41,420	111,635	0.84	0.79–0.90	2.90×10^{-4}
Abdominal aortic aneurysm	6,052	130,888	0.75	0.65–0.87	5.34×10^{-4}
Peripheral vascular disease	13,018	130,888	0.87	0.79–0.96	0.004
Dermatologic					
Psoriasis	5,782	120,136	0.93	0.81–1.07	0.295
Actinic keratosis	31,863	113,000	1.00	0.94–1.06	0.91
Atopic dermatitis	20,684	120,136	0.93	0.86–1.00	0.062
Digestive					
Ulcerative colitis	1,600	114,121	0.96	0.75–1.24	0.768
Diverticulosis	17,826	114,121	0.98	0.91–1.07	0.673
Cholelithiasis	5,371	164,181	0.92	0.79–1.06	0.235
Endocrine/Metabolic					
Graves' disease	561	147,058	1.03	0.88–1.58	0.89
Type 2 diabetes	57,100	95,779	0.88	0.83–0.93	2.50×10^{-4}
Gout	13,865	158,027	0.97	0.89–1.06	0.472
Genitourinary					
Chronic renal failure	14,974	130,019	0.89	0.81–0.97	0.011
Pyelonephritis	619	135,028	0.79	0.51–1.23	0.296
Urinary calculus	12,097	158,884	0.94	0.85–1.03	0.172
Hematopoietic					
Anemia of chronic disease	2,446	134,013	0.91	0.73–1.12	0.358
Neutropenia	1,137	153,778	0.79	0.57–1.10	0.161
Lymphadenitis	2,788	153,778	1.15	0.96–1.37	0.135
Mental disorders					
Schizophrenia	4,010	81,562	0.90	0.76–1.07	0.222
Bipolar	10,919	81,562	0.99	0.89–1.09	0.81
Depression	29,869	81,562	0.96	0.90–1.03	0.257
Musculoskeletal					
Spiral stenosis	11,474	120,114	0.92	0.83–1.02	0.099
Osteoarthritis	10,498	90,509	0.96	0.89–1.04	0.373
Osteoporosis	6,294	157,457	0.89	0.78–1.02	0.002
Neoplasms					
Colon cancer	3,081	112,341	1.03	0.86–1.23	0.783
Cancer of bronchus	2,356	173,071	0.99	0.80–1.21	0.896
Squamous cell carcinoma	2,093	131,987	0.96	0.77–1.20	0.724
Neurological					
Sleep apnea	41,772	102,694	1.06	1.00–1.12	0.063
Migraine	9,330	162,283	1.05	0.94–1.16	0.374
Epilepsy	1,486	145,313	1.04	0.81–1.35	0.735
Respiratory					
Pneumonia	1,430	144,179	1.00	0.77–1.30	0.997
Asthma	12,417	146,715	0.97	0.88–1.07	0.532
Chronic airway obstruction	8,888	146,715	0.88	0.79–0.99	0.032

MVP Publications

 ARTICLE https://doi.org/10.1038/s41467-019-11704-w	 ARTICLES https://doi.org/10.1038/s41588-019-0407-x A catalog of genetic loci associated with kidney function from analyses of a million individuals	 ARTICLES https://doi.org/10.1038/s41588-018-0205-x Corrected: Publisher Correction Genetic analysis of over 1 million people identifies 535 new loci associated with blood pressure traits
Mapping eGFR loci to the renal transcriptome and phenome in the VA Million Veteran Program Jacki  ARTICLES https://doi.org/10.1038/s41588-018-0222-9		 LETTERS https://doi.org/10.1038/s41588-019-0519-3 Genome-wide association analysis of venous thromboembolism identifies new risk loci and genetic overlap with arterial vascular disease
 ARTICLES https://doi.org/10.1038/s41588-018-0303-9 Genetics of blood lipids among ~300,000 multi-ethnic participants of the Million Veteran Program		
 ARTICLES https://doi.org/10.1038/s41588-018-0303-9 Trans-ethnic association study of blood pressure determinants in over 750,000 individuals		 Exome-wide association study of plasma lipids in >300,000 individuals
	 LETTERS https://doi.org/10.1038/s41591-019-0492-5	ARTICLE https://doi.org/10.1038/s41467-019-12576-w OPEN International meta-analysis of PTSD genome-wide association studies identifies sex- and ancestry-specific genetic risk loci
 Genome-wide association study of peripheral artery disease in the Million Veteran Program		  ARTICLES https://doi.org/10.1038/s41588-019-0504-x
Corrected: Author correction; Author correction ARTICLE https://doi.org/10.1038/s41467-019-09480-8 OPEN Genome-wide association study of alcohol consumption and use disorder in 274,424 individuals from multiple populations		Article Published: 29 Jul Target genes, variants, tissues and transcriptional pathways influencing human serum urate levels Genome-wide association study of post-traumatic stress disorder reexperiencing symptoms in >165,000 US veterans

Developing a VA CVD Risk Calculator in the VA Health Care System

Table 3. Harrell C Statistics for Composite ASCVD Events

Model	C statistic (SD)			
	Men		Women	
	White	Black	White	Black
Overall cohort of 1 672 336 veterans				
ASCVD events, No.	54 550	10 575	1154	326
No. at risk	1 314 938	260 225	69 055	28 118
Model 1, 2013 PCE	0.66 (0.004)	0.72 (0.007)	0.78 (0.020)	0.79 (0.036)
Model 2, 2013 PCE with cohort-derived β	0.67 (0.004)	0.72 (0.007)	0.80 (0.018)	0.80 (0.030)
Model 3, 2013 PCE with cohort-derived β and statin therapy	0.67 (0.004)	0.72 (0.007)	0.80 (0.018)	0.80 (0.029)
Subset aged 40-79 y with 1 415 057 veterans				
ASCVD events, No.	48 169	9609	847	285
No. at risk	1 136 161	218 463	44 399	16 034
Model 1, 2013 PCE	0.63 (0.004)	0.68 (0.008)	0.72 (0.022)	0.72 (0.045)
Model 2, 2013 PCE with cohort-derived β	0.64 (0.004)	0.68 (0.008)	0.73 (0.023)	0.73 (0.038)
Model 3, 2013 PCE with cohort-derived β and statin therapy	0.64 (0.004)	0.68 (0.008)	0.73 (0.023)	0.73 (0.036)

Developed at VA specific risk equation for prediction of intermediate and long-term CVD risk using EHR data.

This enables automated estimation of CVD risk in real time using EHR data for use in clinical practice.

Abbreviations: ASCVD, atherosclerotic cardiovascular disease; PCE, Pooled Cohort Equation.

GWAS by PheWAS: Study Design

Genotype Data

- Release 4 AGR + 1000G imputation
- Imputation info score >0.3
Minor allele count > 20

PheWAS

AFR

AMR

EUR

Multi-population
Meta

Summary of GWS Associations

Independent Lead SNPs

Heritability

Estimated compute time was 8 years on the available computing infrastructure

Genetic Correlations

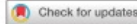
Pleiotropic Associations

Heterogeneity Analysis

- 1,866 Diagnosis Codes (Phecodes)
132 question from core MVP enrollment questionnaire
- 6 vital measurement

Using SAIGE - a generalized mixed model association test that uses the saddle point approximation to

Ancestries: African (AFR), Admixed Americans (AMR), East Asian (EAS), European (EUR)



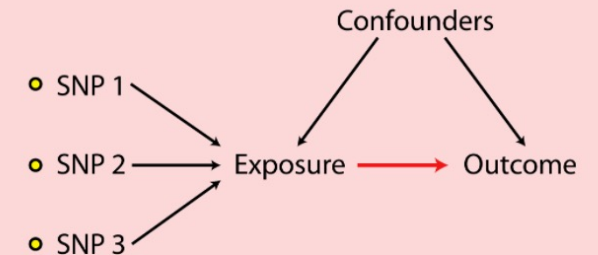
Actionable druggable genome-wide Mendelian randomization identifies repurposing opportunities for COVID-19

Liam Gaziano^{1,2}, Claudia Giambartolomei^{3,4}, Alexandre C. Pereira^{5,6}, Anna Gaulton⁷, Daniel C. Posner⁸, Sonja A. Swanson⁸, Yuk-Lam Ho¹, Sudha K. Iyengar^{9,10}, Nicole M. Kosik¹, Marijana Vujkovic^{11,12}, David R. Gagnon^{1,13}, A. Patrícia Bento¹⁷, Inigo Barrio-Hernandez¹⁴, Lars Rönnblom¹⁵, Niklas Hagberg¹⁵, Christian Lundtoft¹⁵, Claudia Langenberg^{16,17}, Maik Pietzner¹⁷, Dennis Valentine^{18,19}, Stefano Gustincich³, Gian Gaetano Tartaglia³, Elias Allara², Praveen Surendran^{2,20,21,22}, Stephen Burgess^{2,23}, Jing Hua Zhao², James E. Peters^{21,24}, Bram P. Prins^{2,21}, Emanuele Di Angelantonio^{2,20,21,25,26}, Poornima Devineni¹, Yunling Shi¹, Kristine E. Lynch^{27,28}, Scott L. DuVall^{27,28}, Helene Garcon¹, Lauren O. Thomann¹, Jin J. Zhou^{29,30}, Bryan R. Gorman¹, Jennifer E. Huffman³¹, Christopher J. O'Donnell^{32,33}, Philip S. Tsao^{34,35}, Jean C. Beckham^{36,37}, Saiju Pyarajan¹, Sumitra Muralidhar³⁸, Grant D. Huang³⁸, Rachel Ramoni³⁸, Pedro Beltrao¹⁴, John Danesh^{2,20,21,25,26}, Adriana M. Hung^{39,40}, Kyong-Mi Chang^{12,41}, Yan V. Sun^{42,43}, Jacob Joseph^{1,44}, Andrew R. Leach⁷, Todd L. Edwards^{45,46}, Kelly Cho^{1,47}, J. Michael Gaziano^{1,47}, Adam S. Butterworth^{2,20,21,25,26} ✉, Juan P. Casas^{1,47} ✉ and VA Million Veteran Program COVID-19 Science Initiative*

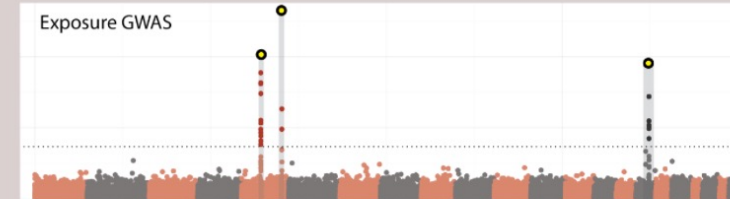
Drug repurposing provides a rapid approach to meet the urgent need for therapeutics to address COVID-19. To identify therapeutic targets relevant to COVID-19, we conducted Mendelian randomization analyses, deriving genetic instruments based on transcriptomic and proteomic data for 1,263 actionable proteins that are targeted by approved drugs or in clinical phase of drug development. Using summary statistics from the Host Genetics Initiative and the Million Veteran Program, we studied 7,554 patients hospitalized with COVID-19 and >1 million controls. We found significant Mendelian randomization results for three proteins (ACE2, $P = 1.6 \times 10^{-6}$; IFNAR2, $P = 9.8 \times 10^{-11}$ and IL10RB, $P = 2.3 \times 10^{-14}$) using *cis*-expression quantitative trait loci genetic instruments that also had strong evidence for colocalization with COVID-19 hospitalization. To disentangle the shared expression quantitative trait loci signal for *IL10RB* and *IFNAR2*, we conducted genome-wide association scans and pathway enrichment analysis, which suggested that *IFNAR2* is more likely to play a role in COVID-19 hospitalization. Our findings prioritize trials of drugs targeting IFNAR2 and ACE2 for early management of COVID-19.

- Perform two-sample Mendelian randomization (MR) for 1,263 “actionable” proteins on COVID-19 hospitalization

Objective: Infer the causal effect of the exposure on the outcome



1.

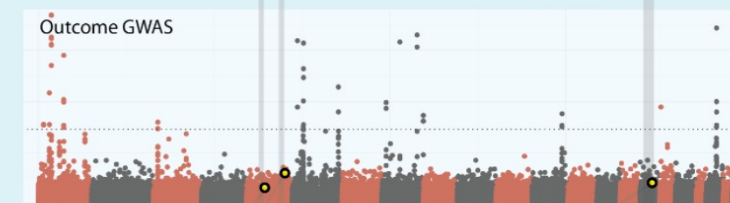


Description

Define instruments: Obtain SNPs that are GWAS significant for the exposure. Ensure that they are independent.

Instruments can be defined from a variety of different sources.

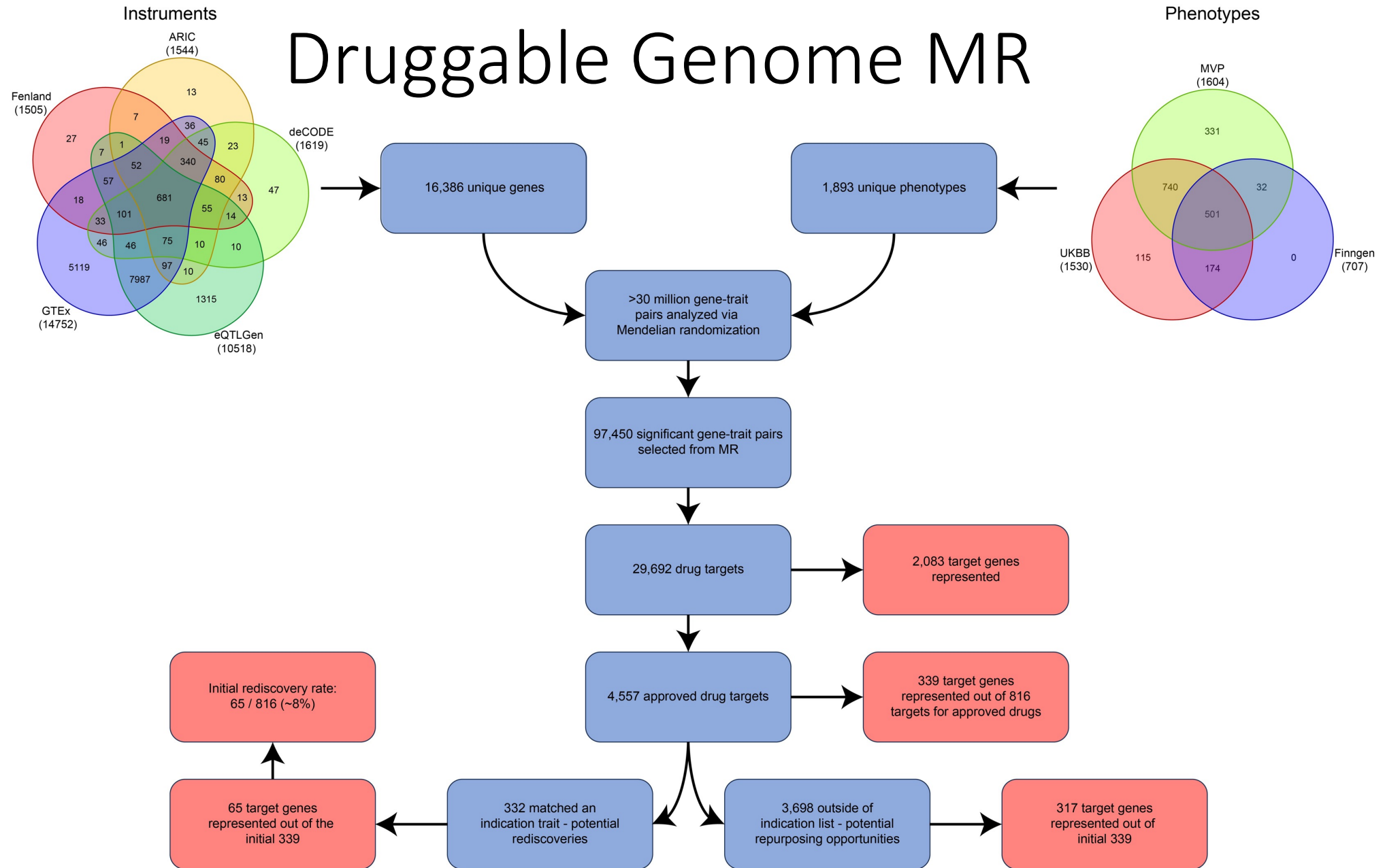
2.



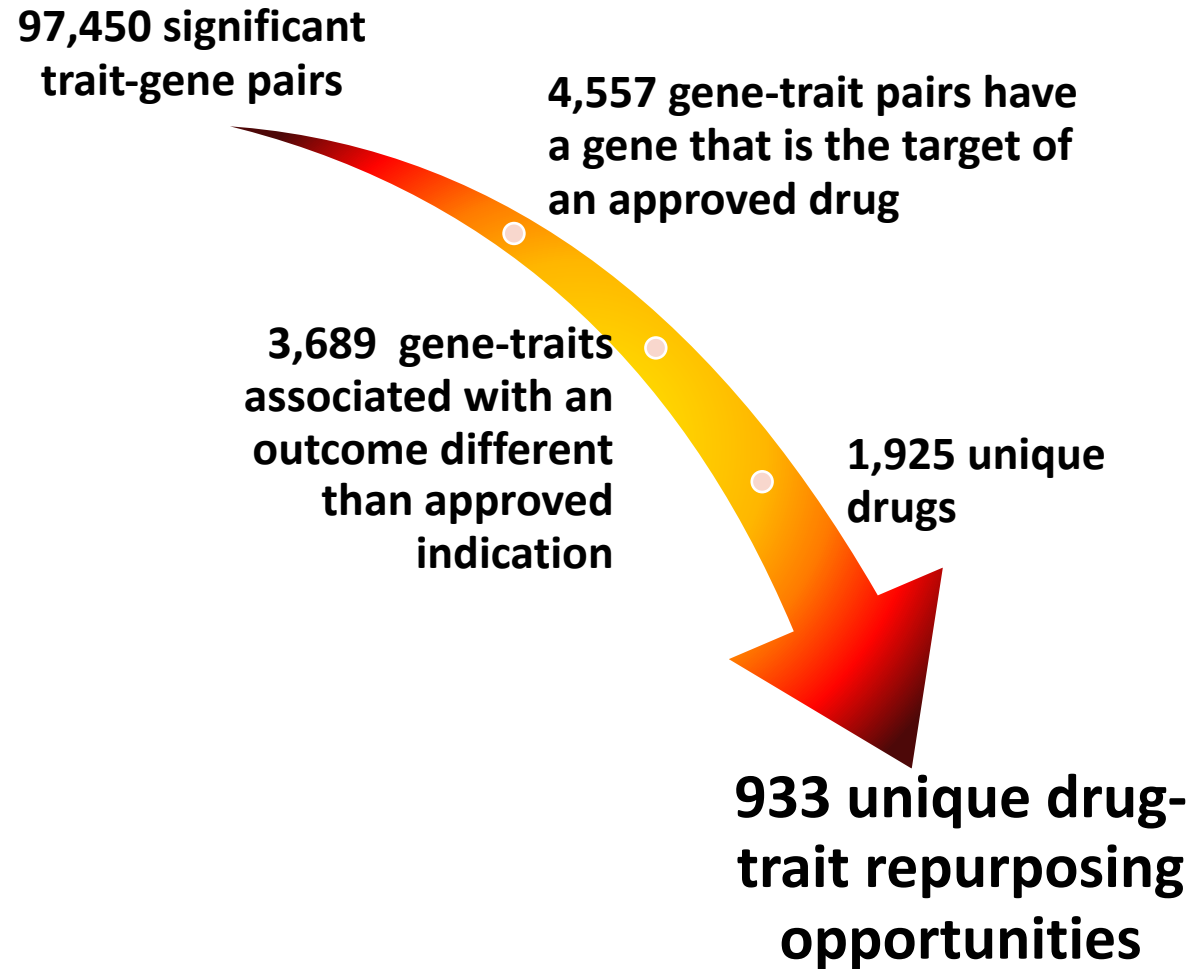
Get effects on outcome: Extract the instrument SNPs from the outcome GWAS. If they are not available, use LD proxies instead.

MR Base contains a large database of entire GWAS summary statistics.

Druggable Genome MR



The repurposing landscape:



Familial Hyperlipidemia Return of Genetic Results in MVP

Genomic Profiling of CV Risk: Precision Prevention and Treatment



Planned: Identify **monogenic** carriers

- LDLR, APOB, PCSK9, all <0.5%
- Early CVD surveillance
- Intense LDL therapy
- Cascade screening

Future: Profiling for high **polygenic** risk

- Early lifestyle and statin Rx
- Lower threshold for CVD screening
- Guideline-based

Is Genetics Ready for Prime Time?

- Cancer Genetics
- Monogenic disease
- Pharmacogenomics
- Prediction/PRS
- Drug Development
- Therapeutics
- Others

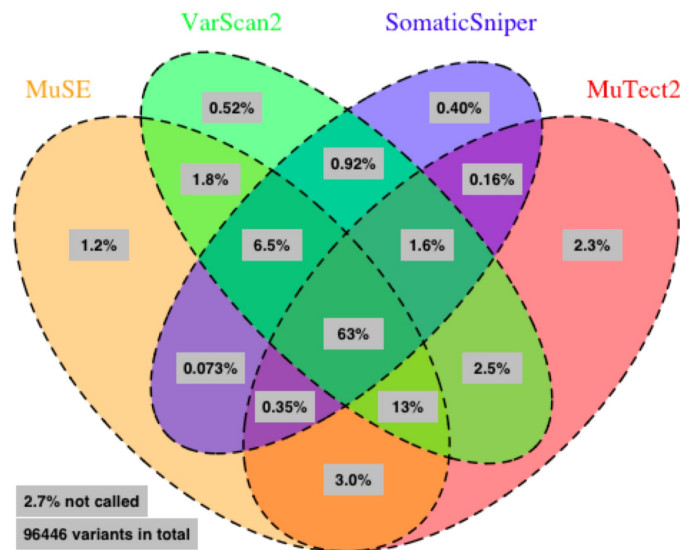
Challenges for Data Analytics

- Structural
 - Maintaining cohorts: costs and priorities
 - Diversity
- Analytic
 - Omics assessment
 - Phenotype and exposome assessment
 - Multidimensional data
 - Analytic and computation challenges
- Cross talk between cohorts

Omics Challenges

- Scale and Cost
- Evolution of Technology
- Consistency over time
- Meta analyzing
- Measurement error

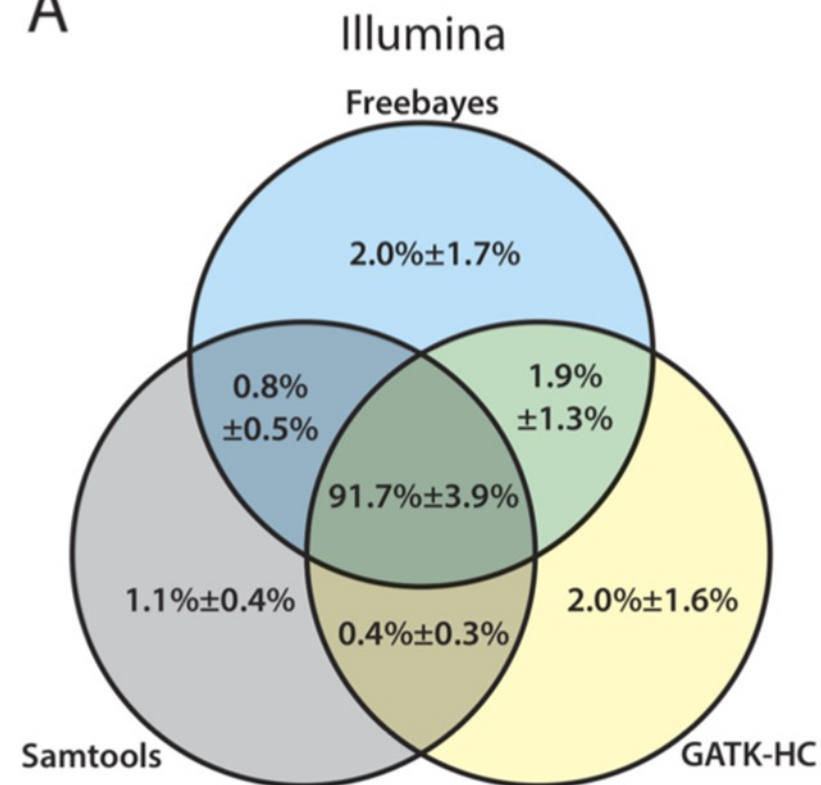
System 3: Data Harmonization System To Analyze all of the Submitted Data with a Common Pipelines



- MuSE (MD Anderson)
- VarScan2 (Washington Univ.)
- SomaticSniper (Washington Univ.)
- MuTect2 (Broad Institute)

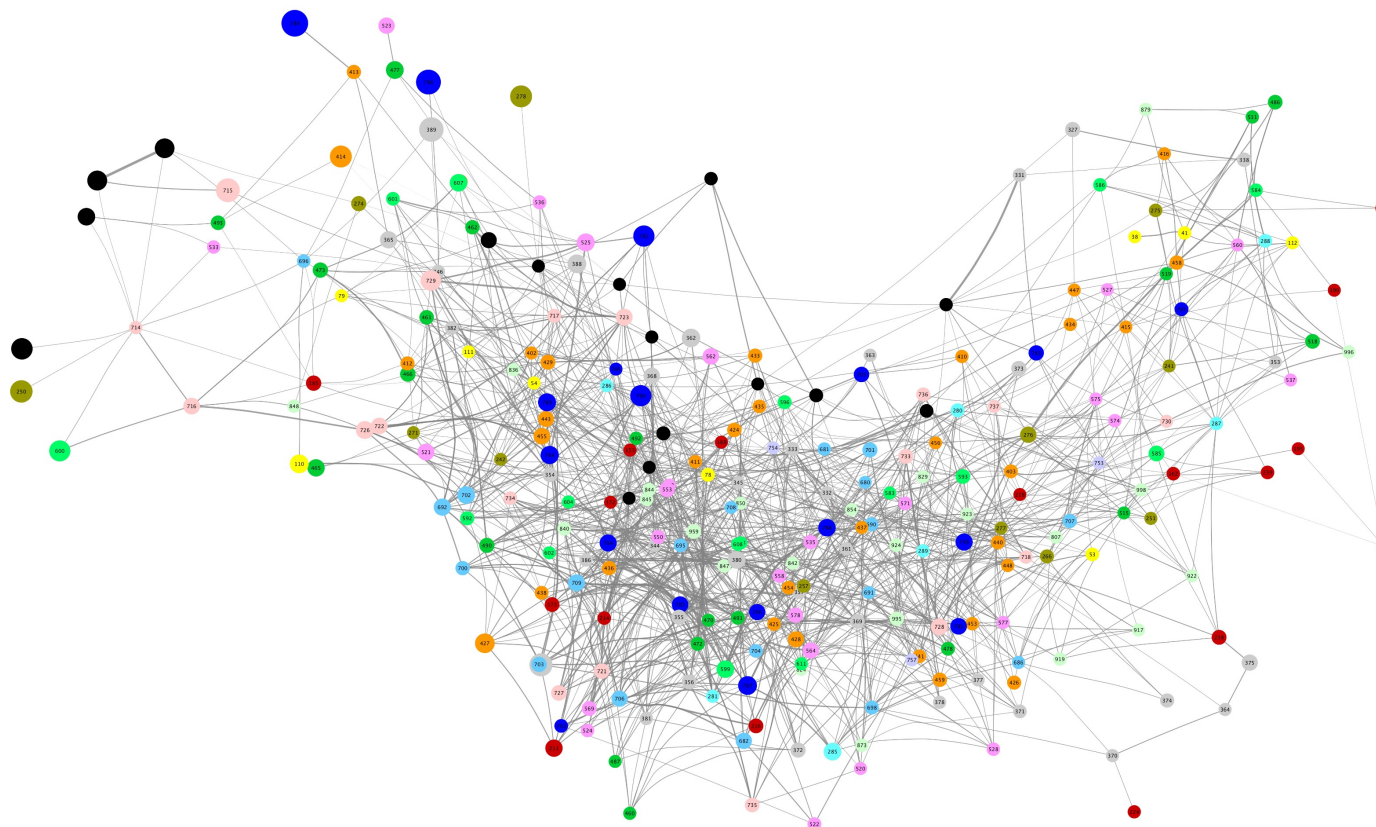
Source: Zhenyu Zhang, et. al. and the GDC Project Team, Uniform Genomic Data Analysis in the NCI Genomic Data Commons, to appear.

A



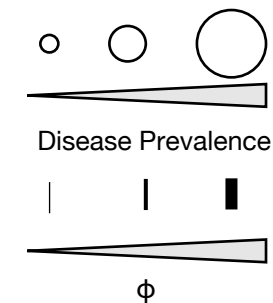
Venn diagrams summarizing called variants by different callers. The mean percentage with standard deviation of confidence variant calls with equal to or higher than the quality score threshold of 20 are represented for (A) Illumina data sets and

VA Phenotypic Network Map



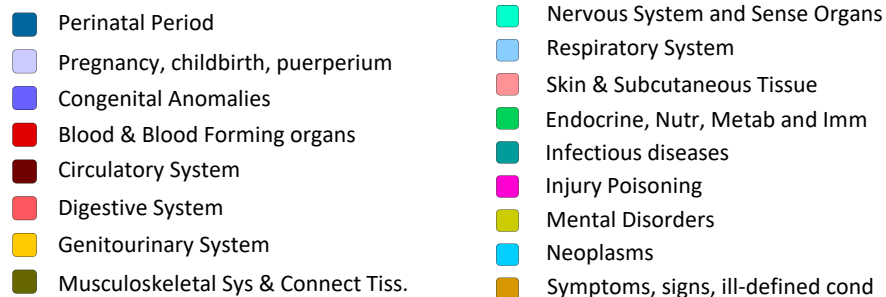
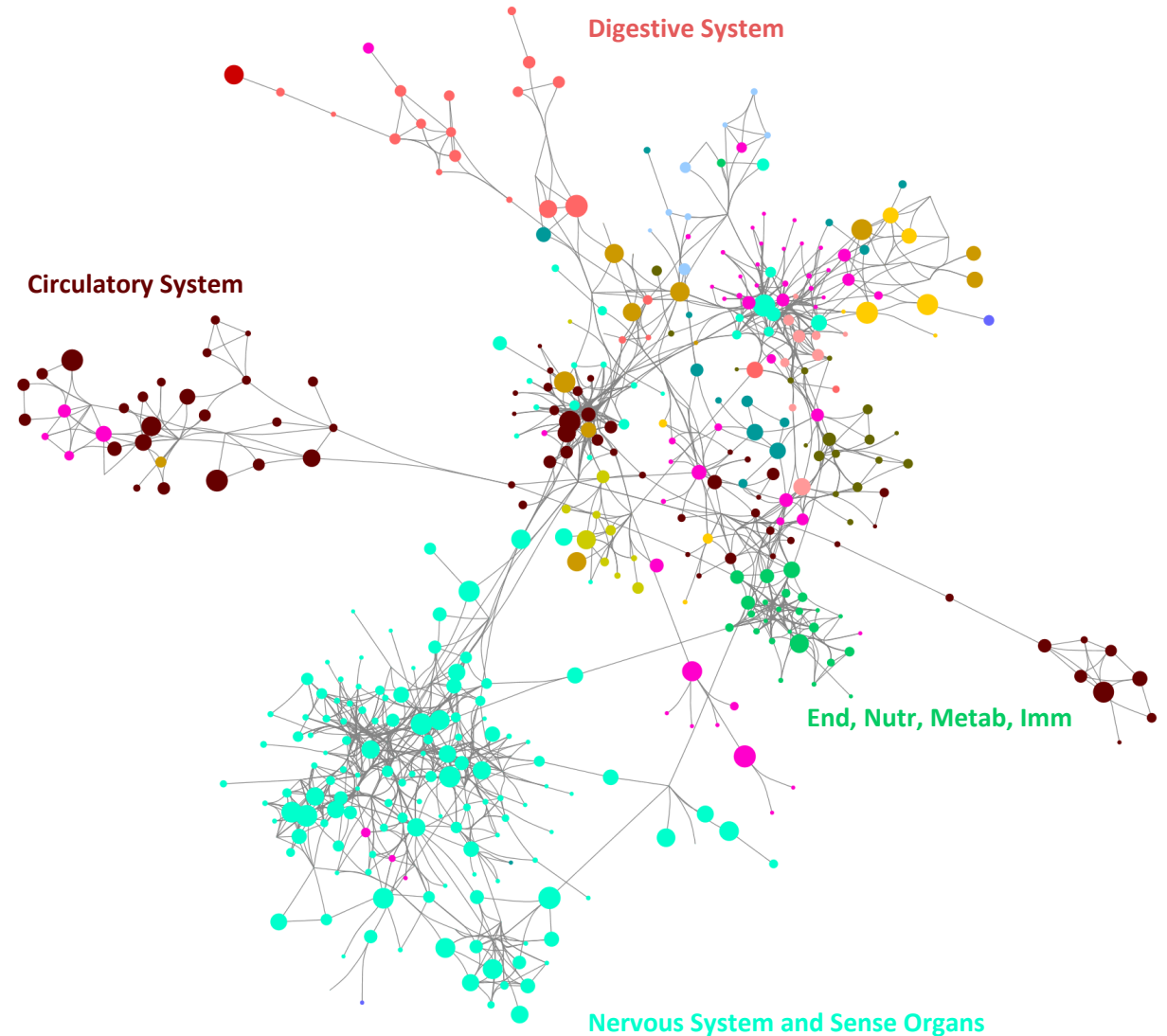
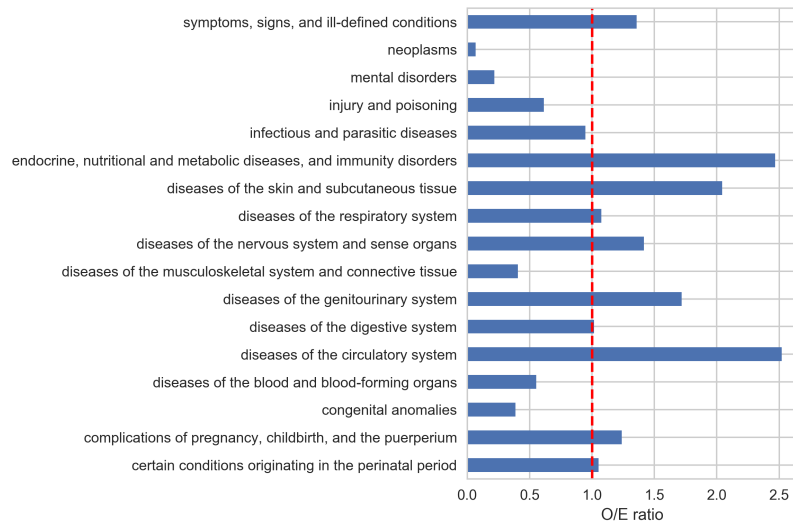
Disease Category

- | | | |
|--------------------------------|--|---|
| Congenital Anomalies | Musculoskeletal System and Connective Tissue | Infections and Parasitic |
| Blood and Blood Forming Organs | Nervous System and Sense Organs | Injury and Poisoning |
| Circulatory System | Respiratory System | Mental Disorders |
| Digestive System | Skin and Subcutaneous Tissue | Neoplasms |
| Genitourinary System | Endocrine, Nutritional, Metabolic and Immunity | Symptoms, Signs, and ill-defined conditions |



Community #3: Circulatory System

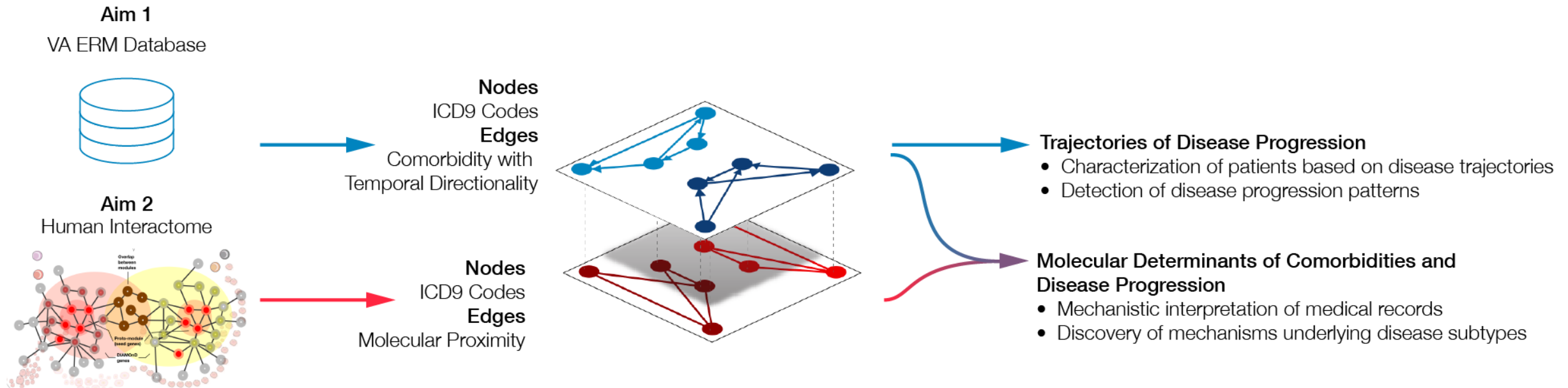
- 830 codes
- Node size: Prevalence



* Links were removed for visualization purposes

Applying Network Medicine in the VA

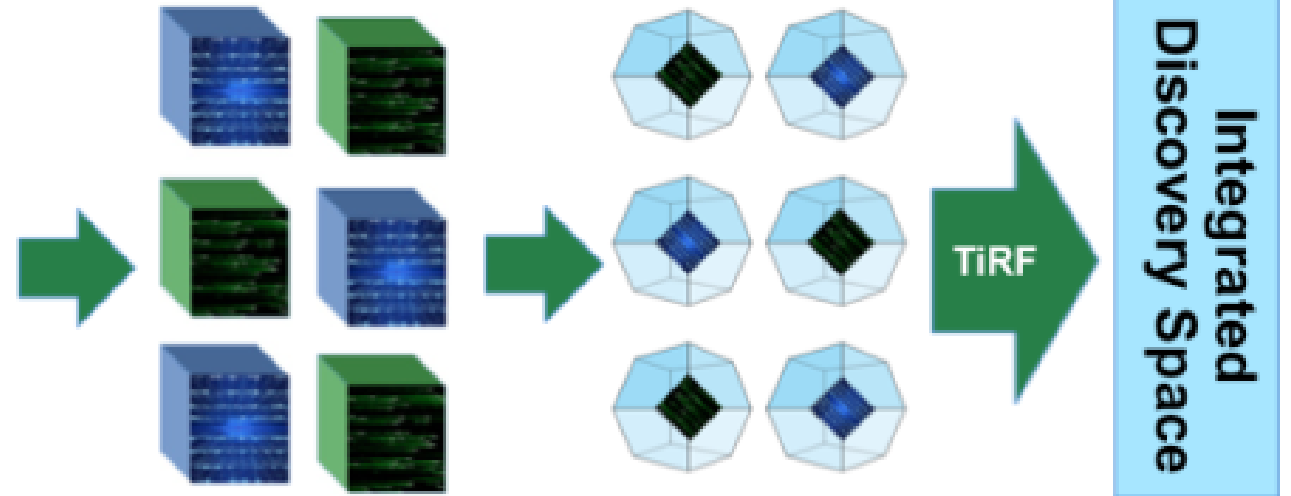
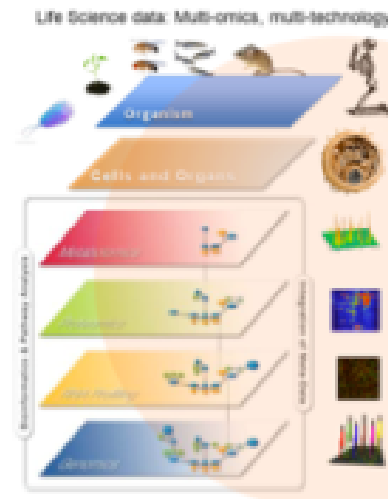
We are applying Network Medicine tools for evaluating disease trajectories and molecular determinants of comorbidities and disease progression



Computing Challenges

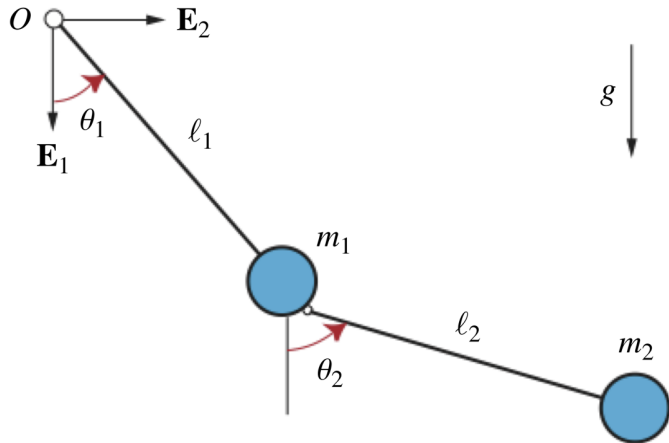
- Capacity
- Cost
- Expertise
- Data movement
- Security
- Access
- Federation
- AI /ML

Discovery: Matrices → Cubes → Polytopes

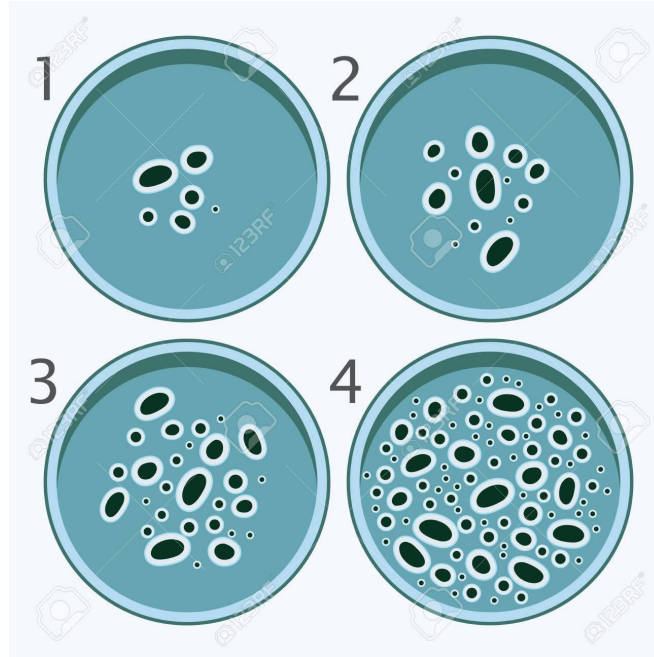


AI Challenges

Double pendulum problem



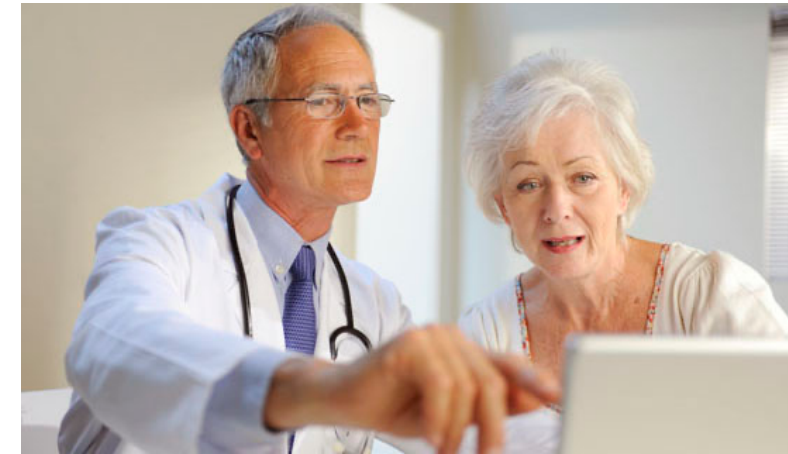
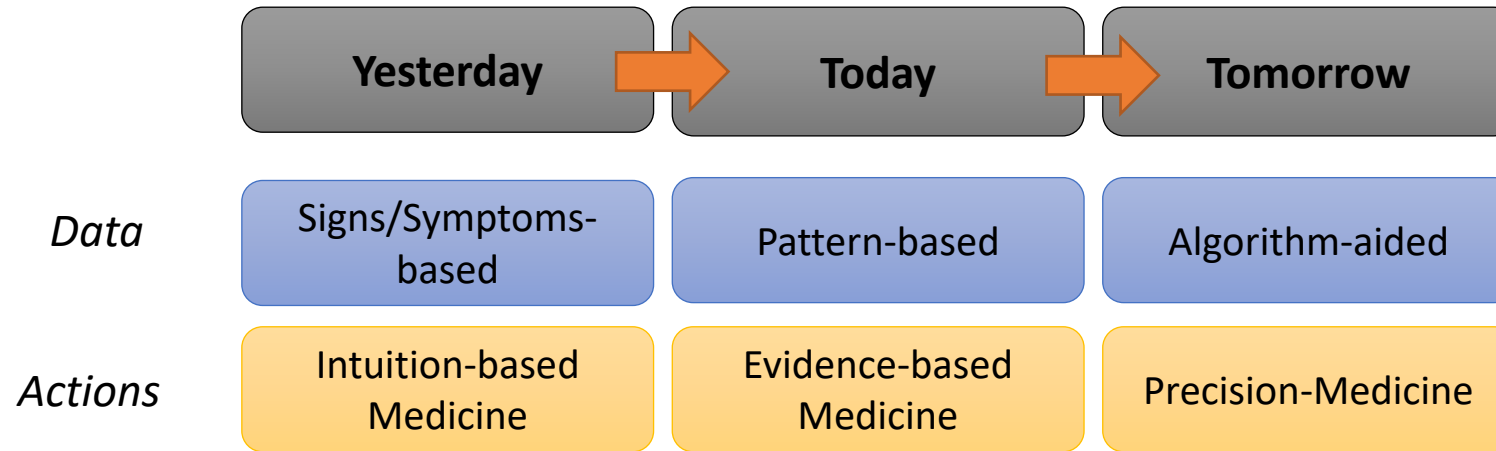
Single organism growing in a culture



Correlated images



The Promise of big health data: To improve health, prevent and treat disease by empowering providers, patients and health systems with better access to and use of health data.



Inside the GE Healthcare-designed "Capacity Command Center" at the John Hopkins Hospital.
Photograph by Ryan Donnell for Fortune



What are the next steps for the use of big data, genetics and other omics in health care systems.

- Clinical data analytics
- Quality improvement
- Research
 - Genetics other omics discovery
 - Prediction, PRS, etc.
- Comparative effectiveness
- Cost efficacy
- Business operations

Requires investment in computing and data infrastructure and unprecedented collaboration between clinicians, researchers, operations and even patients



"I'm participating in the Million Veteran Program so that I can do my part to help future generations of not just Veterans, but everyone who can benefit from this research."



"When I was young, the service gave me a reason. Today, my reason is to help answer those questions yet to be asked....."



"I have always known someone in the family with Diabetes or Hypertension. I eagerly volunteered to participate in MVP so I can help medical researchers better understand how genes influence diseases. One blood draw is all it took... yet the potential to contribute to scientific discoveries is enormous!"



"I believe that the data collected from me and other Veterans in the Million Veterans Program will someday provide better ways to diagnose and treat patients. I volunteered to participate in the MVP because I want to do my part to make this a reality one day."



Thank You Discussion

