

# Applications of Machine Learning in EHR for HIV Subphenotyping

@

## COERE 2025 Annual Methods Symposium

Sandra Safo ([ssafo@umn.edu](mailto:ssafo@umn.edu); [www.sandraesafo.com](http://www.sandraesafo.com))

Division of Biostatistics and Health Data Science

April 10, 2025

Financial Disclosure: None



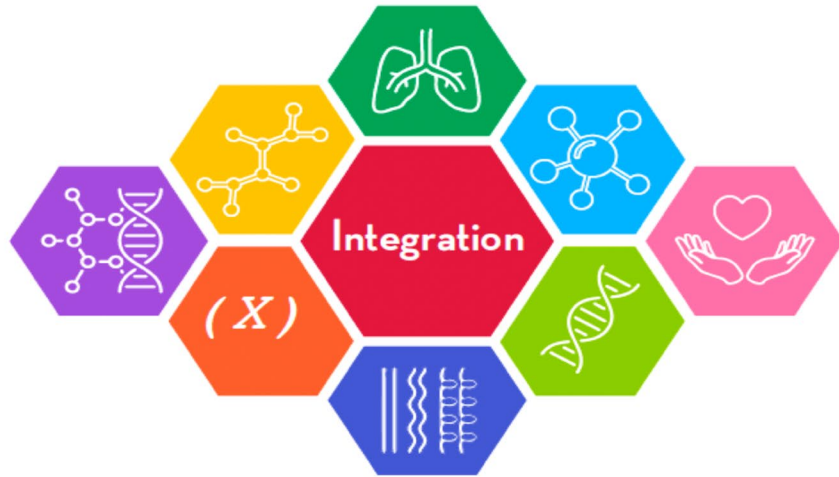
UNIVERSITY OF MINNESOTA

**Driven to Discover®**

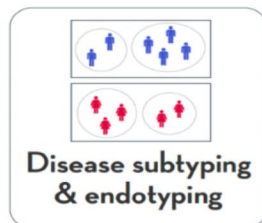
# Research Overview



Machine and  
Statistical Learning



Disease prediction



Disease subtyping  
& endotyping



Greater insight  
into disease  
mechanism



Potential  
therapeutic  
targets



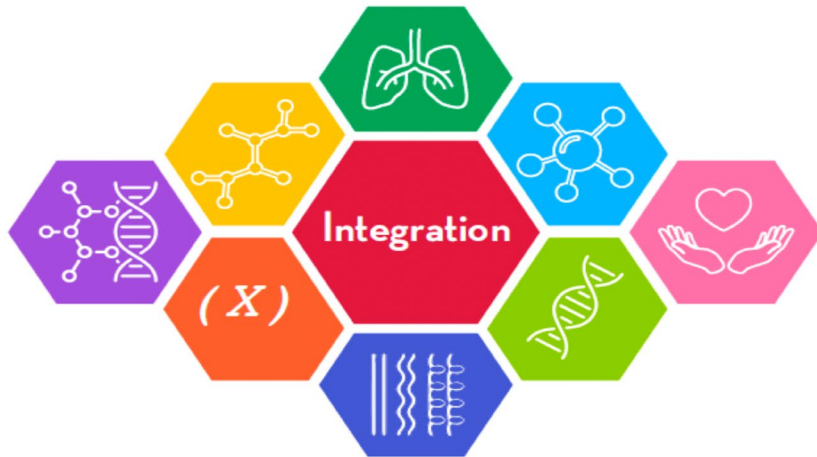
Dissemination



# Research Overview



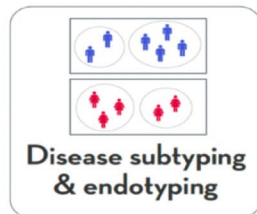
Machine and  
Statistical Learning



## ***From Big Data to Health Discovery***



Disease prediction



Disease subtyping  
& endotyping



Greater insight  
into disease  
mechanism



Potential  
therapeutic  
targets



Dissemination



# Roadmap

Summary and  
Key Takeaways

How can you get started?

Challenges and Considerations

A Gentle Intro to ML  
And EHR

ML in Action



[This Photo](#) by Unknown Author is licensed under [CC BY-SA](#)



# Machine Learning

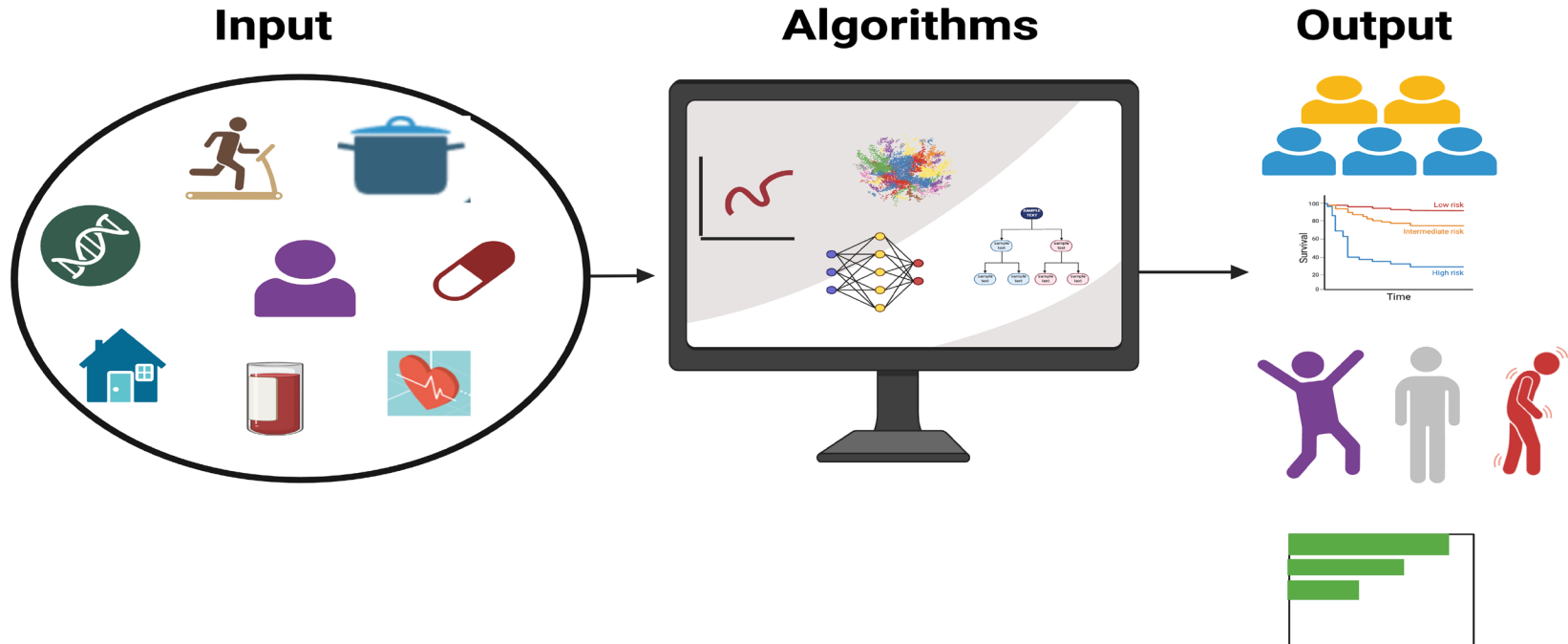


Figure generated with BioRender



# Supervised ML

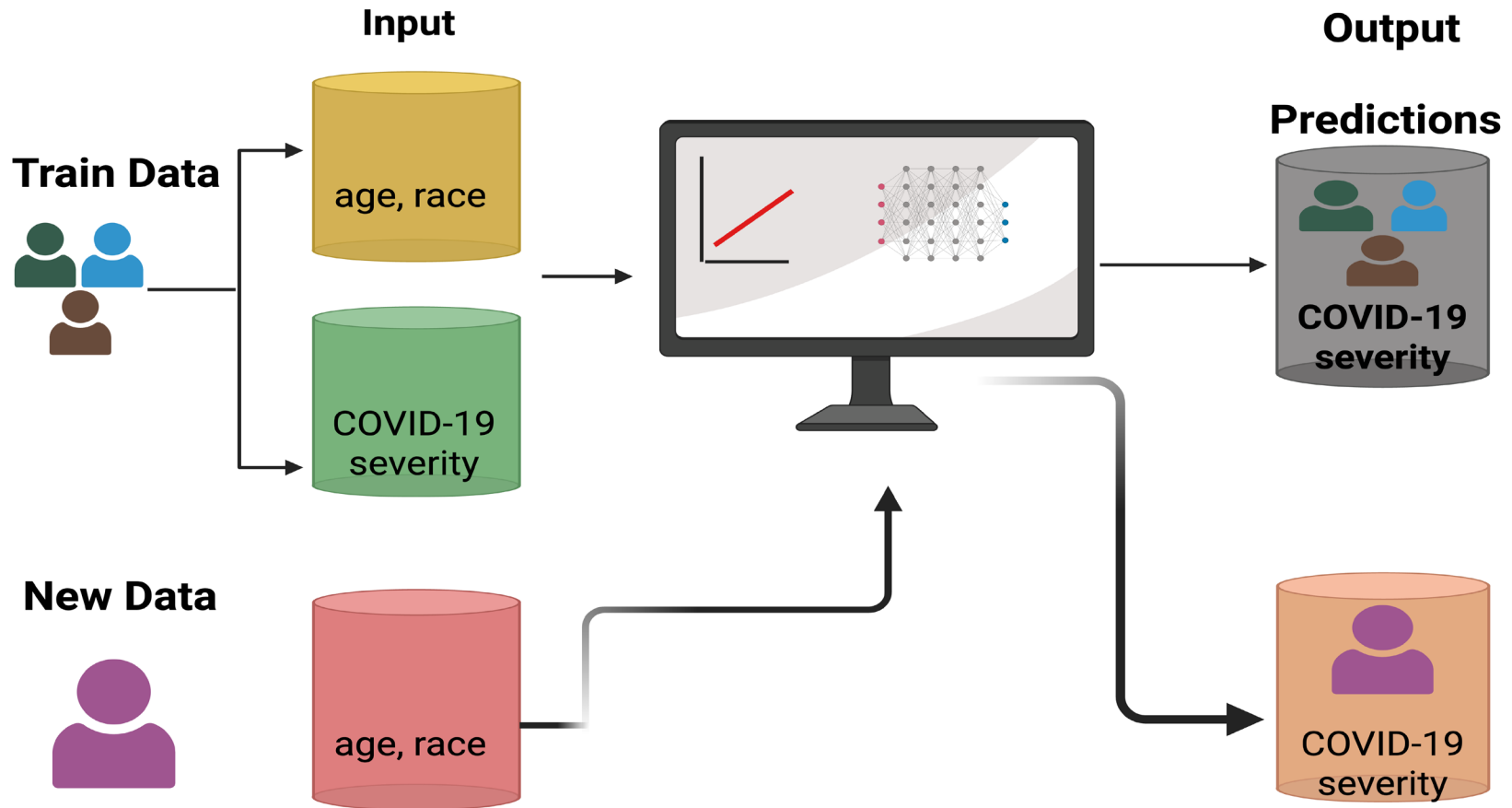


Figure generated with BioRender

# Unsupervised ML

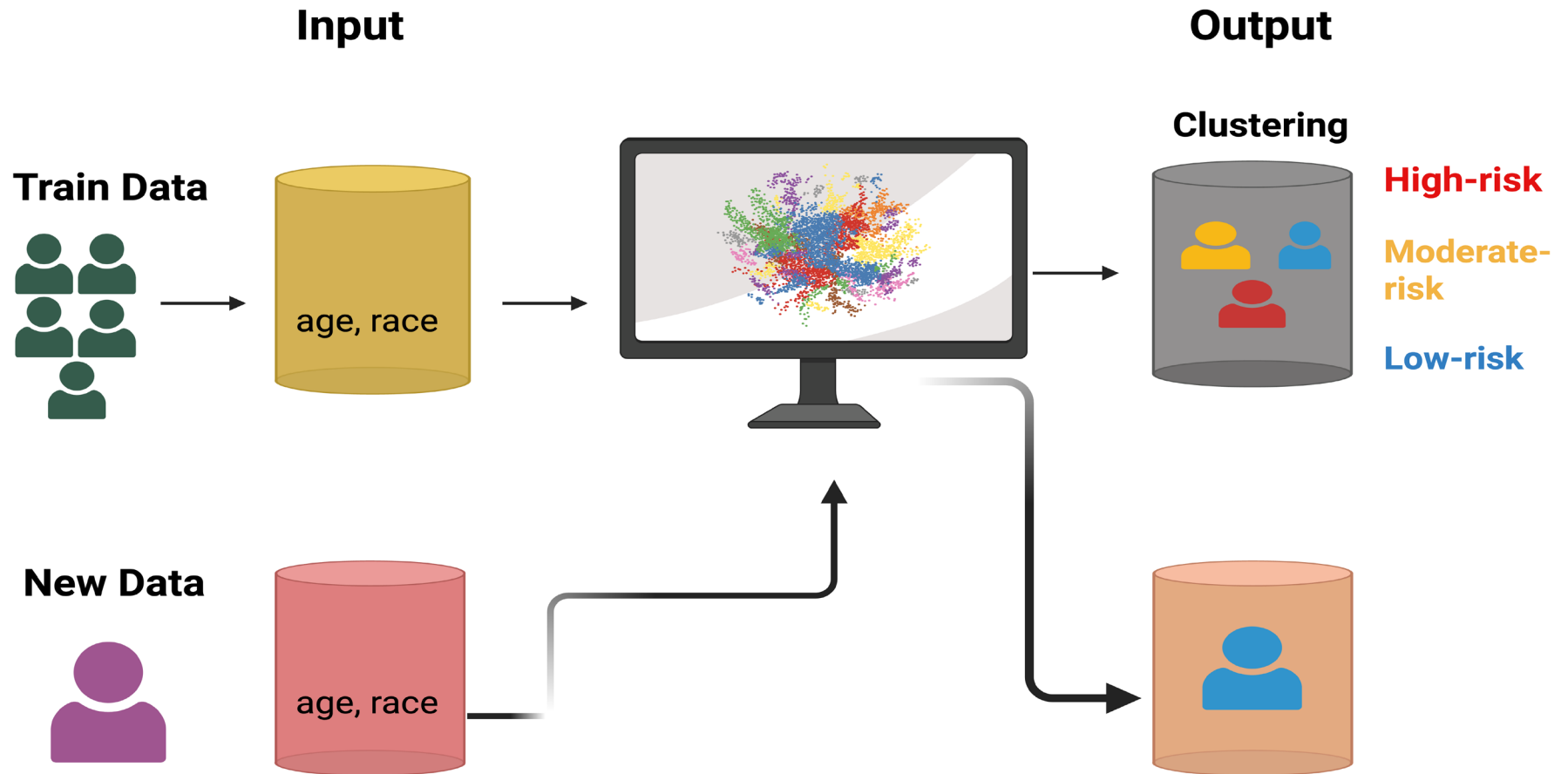


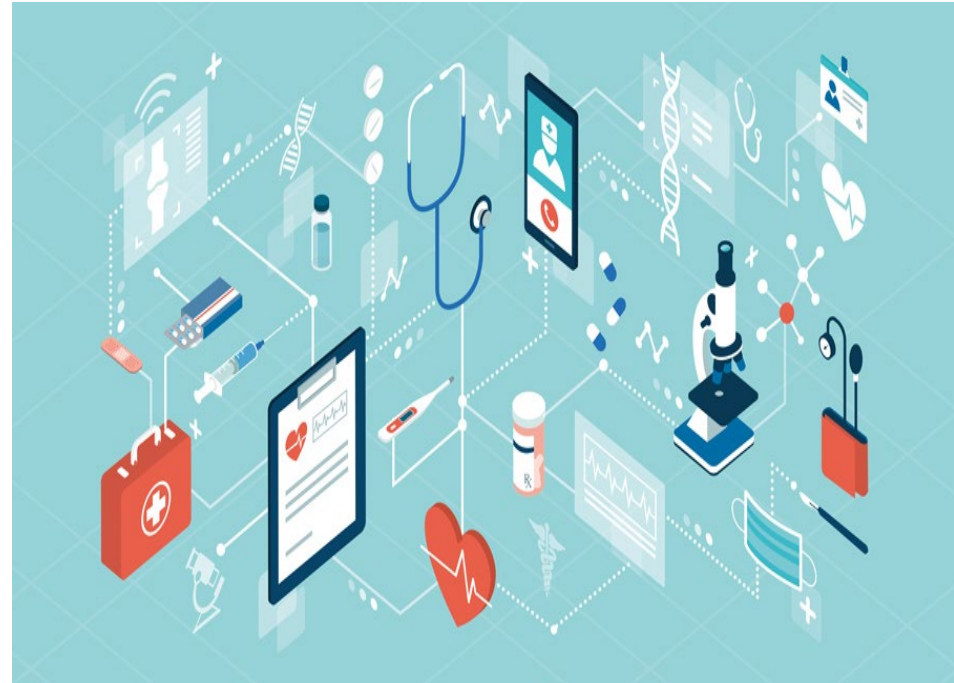
Figure generated with BioRender



# Why care about ML?

## Some Challenges in Health Care

- Diagnosing diseases early
- Predicting patient outcomes (e.g., risk of severe COVID-19, increased hospitalizations, risk of Long COVID)
- Personalizing treatments based on patient data
- Patient selection for clinical trials



[This Photo](#) by Unknown Author is licensed under [CC BY-NC-ND](#)





# Why care about ML?

## Opportunities with ML

- Better diagnostic tools
- Improved decision support
- Data-driven research eliminating guestimates
- Help uncover hidden patterns in large and complex data (e.g., **EHR**)



[This Photo](#) by Unknown Author is licensed under [CC BY-SA](#)



# Challenges and Opportunities of EHR



**EHR- digitized version of a patient's chart and clinical records over time**



**Structured and Unstructured Data**

Demographics,  
Medications, vital signs, past medical history,  
Comorbidities,  
laboratory data,  
doctor's note etc



**Challenges**

Large, multidimensional, and unstructured  
Missingness  
Data quality (e.g., misclassification, measurement error, selection bias)



**Opportunities**

Can be cost effective to maintain  
Data retrieval is quick  
Offers opportunities to answer research questions and improve patient outcomes



# **Machine Learning in Action: COVID-19 Subphenotyping among Persons Living with HIV (PLWH)**





---

Persons living with HIV (PLWH) are disproportionately impacted by COVID-19

---

Older age, male sex, a history of smoking and comorbidities, is associated with risk of severe clinical outcomes or death

---

PLWH are more likely to have multiple comorbidities

---

Impact of factors, e.g. clinical comorbidities, predictive of severe COVID-19 among PLWH varies by race/ethnicity

---

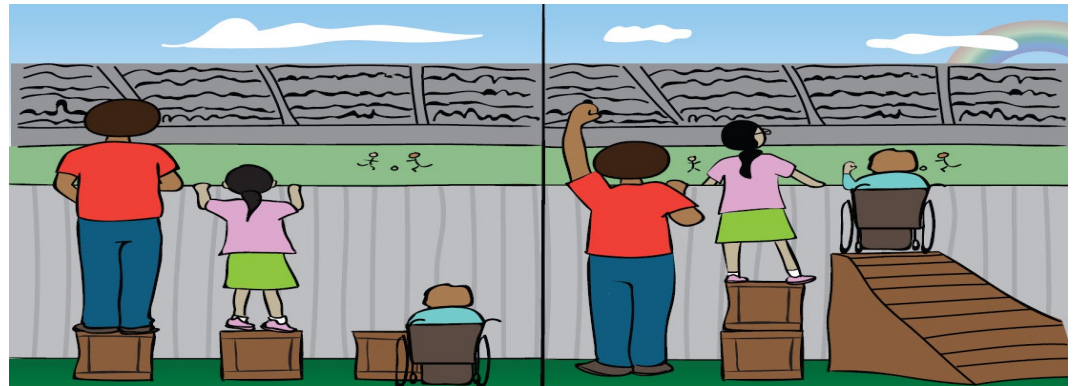
Limited work exists that comprehensively investigate multiple data modalities including demographics, social determinants of health (SDoH), comorbidities, to better understand the characteristics of PLWH who tend to have worse outcomes.

## **Background and Rationale of Study**





- To **discover** and **validate** clusters with varying risks (COVID-19 severity, death, hospitalizations, Long COVID)
- To investigate the variables characterizing these clusters
- Our findings could
  - Facilitate targeted policy interventions
  - Enable a deeper understanding of health disparities in PLWH and promote health equity



# N3C Enclave: One of the largest public HIPAA- limited data set in US history

10/10/24



National  
COVID  
Cohort  
Collaborative



**Sites: 84**



**Persons: 22.8 million**



**COVID+ Cases: 8,914,402**



**# of Rows: 33.9 billion**



**Clinical Observations:  
3.3 billion**



**Lab Results: 16.3 billion**



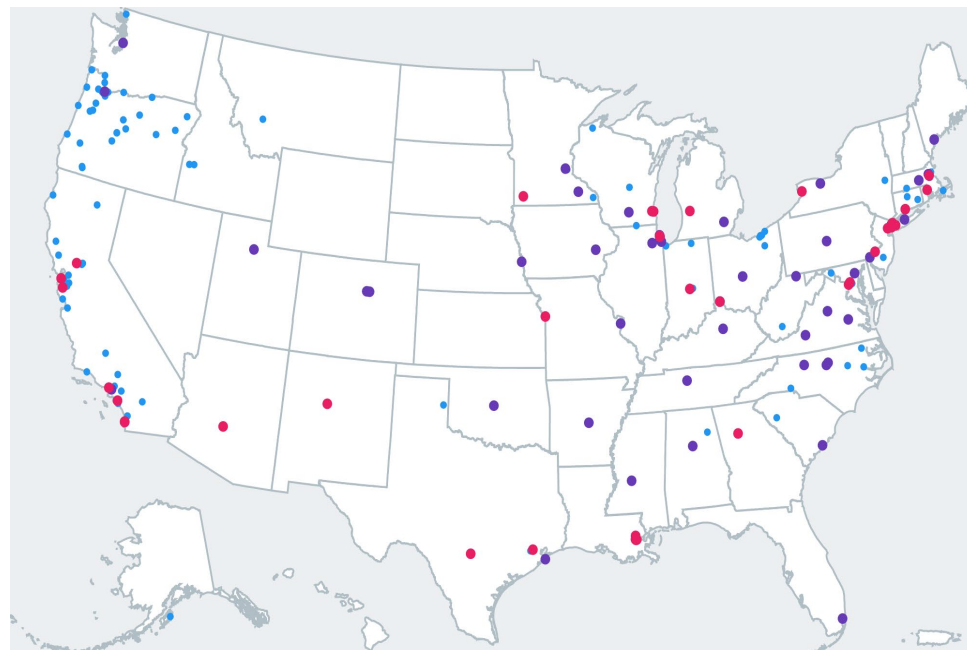
**Medication Records:  
5.3 billion**



**Procedures: 1.2 billion**



**Visits: 2.0 billion**



# Data



Data from N3C (National COVID Cohort Collaborative)

Spans multiple states and hospitals  
January 1, 2020 to November 2, 2023



Exclusions

Age < 18  
Individuals taking medications for HIV prevention or treatment of chronic hepatitis B infection



**Outcome:** COVID-19 Severity

**Not Severe:** Asymptomatic to symptomatic assistance needing hospitalization

**Severe:** Hospitalized with invasive mechanical ventilation, extracorporeal membrane oxygenation [ECMO], discharge to hospice or death

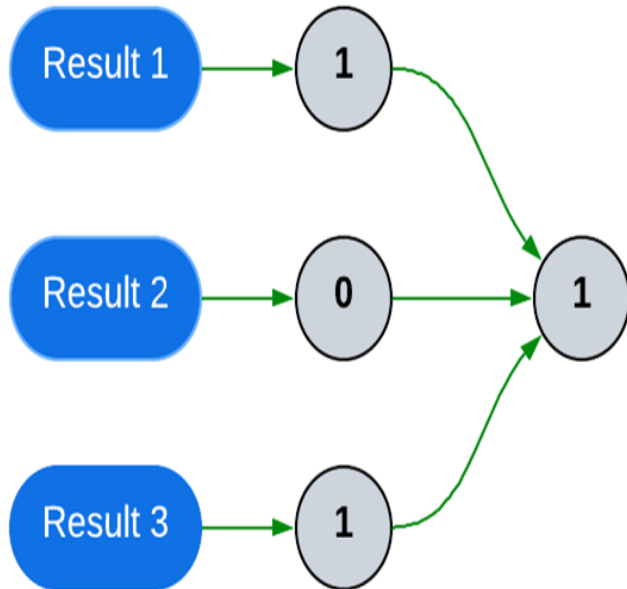


Input Variables: demographics, comorbidities, lab measurements, social determinants of health (SDoH)





# Methods



## Data splitting (Focus on PLWH)

- Discovery set: to detect the clusters
- Testing set: to reproduce the clusters

## Variable selection

- Identify which input variables discriminate COVID-19 severity group
- Use these variables to cluster individuals
- Incorporate resampling techniques for statistical rigor and address missingness

## Cluster Detection

## Reproduce Clusters





# Variable Ranking via Ensemble of Methods

Variable Importance by Methods

Variable Names	Lasso	Elastic Net	Random Forest	XG Boost	Light GBM	Overall
Bilirubin total (mg/dL)	39.4	39.2	40	40	40	39.7
Albumin (g/dL)	37.8	37.8	39	39	39	38.5
ALT/SGPT (IU/L)	39.6	39.8	38	38	37	38.5
CHF	37	37	33.8	37	38	36.6
Hemoglobin (g/dL)	35	34.8	36.8	35.2	31.2	34.6
Temperature (degrees C)	36	36	29.4	34.4	33.4	33.8
Glucose (mg/dL)	31.4	31.6	33.4	35.4	34.4	33.2
Renal Disease	32.6	32.6	26.2	32.4	36	32
Respiratory rate (BPM)	27.2	27.6	26.8	29.2	32.8	28.7
Creatinine (mg/dL)	27.2	27.4	27.8	31.6	23	27.4
Current Smoker	34	34.4	10.4	26	25	26
Stroke	29.4	28.8	16.2	21.8	26	24.4
Diastolic blood pressure	20	20.2	32.6	22.2	19.8	23
Systolic blood pressure (mmHg)	21.6	20.6	32	19.4	17.6	22.2
AST/SGPT (IU/L)	18	21.6	21	20.2	29.2	22
BMI before COVID	17.6	19.4	21.2	28.4	21.8	21.7
Lymphocytes (x10E3/uL)	10.8	9.6	21.6	31.8	32.4	21.2
Race Ethnicity: Hispanic	31.4	31.4	9.8	14.4	17.2	20.8
Sodium (mmol/dL)	16.2	15.4	30.6	20.2	20.4	20.6
Chloride (mmol/dL)	11.4	9	30.4	24.6	25.8	20.2
Age (years)	10.8	8.2	36.2	24.8	17.6	19.5
Myocardial Infarction	26	24.2	15	14.2	17.8	19.4
Platelet count (x10E3/uL)	15.2	13.8	21.8	26.8	17.4	19
cci_dmcx	23.8	23.6	17.6	10.6	17.2	18.6
Mild Liver Disease	24.4	23.6	6.8	17.8	19.6	18.4
Pulmonary	25.6	25.8	2.4	11	24.6	17.9
White blood cell count (x10E3/uL)	11	8.4	12.8	25.6	25.6	16.7
Metastatic Cancer	20.4	19.6	15.6	8.2	19.6	16.7
Paralysis	21.4	20.6	8.6	7.6	17.2	15.1
BUN/Creatinine ratio	11.2	8.8	23.4	14	17.2	14.9
Gender: male	16	16.8	6.4	7.6	26.6	14.7
Peripheral Vascular Disease	15.8	14.2	14.2	10.4	17.2	14.4
Cancer	19.6	18.4	11	3.4	17.4	14
Potassium (mmol/L)	10.8	9.4	12.6	15.8	20.6	13.8
Diabetes	12.4	11.2	21	2.4	17.6	12.9
Severe Liver Disease	10.8	8.2	20.6	5	17.2	12.4
Neutrophils (x10E3/uL)	11.4	9	2.6	20.2	17.2	12.1
SpO2	12.8	10.2	4.6	8.8	18.8	11
Peptic Ulcer Disease	12.4	11.4	7.8	3.4	17.2	10.4
Rheumatic Disease	13	12.6	2	2.4	17.2	9.4

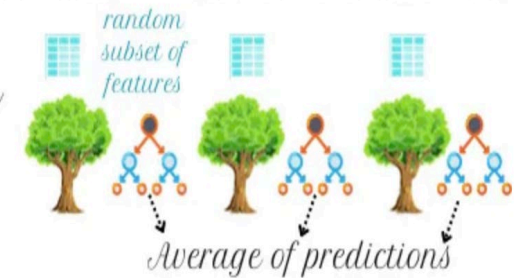
Mean Score for Methods

LASSO (Least absolute shrinkage and selection operator), Elastic Net, Random Forest, XG Boost, Light GBM

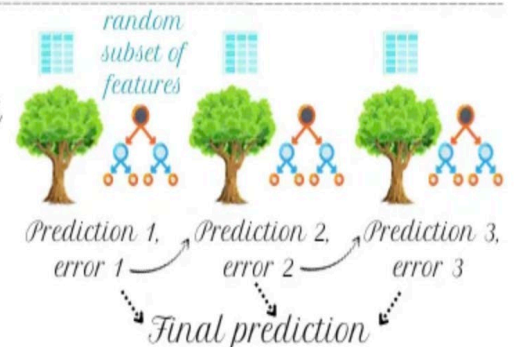
Decision Tree



Random Forest



XGBoost



Source: Medium



# Baseline Characteristics

	All N3C (n=5622302)	Persons without HIV (n=5547210)	Persons with HIV (n=75092)
<b>COVID-19</b>			
Not Severe	5527166 (98%)	5453251 (98%)	73,915 (98%)
Severe	95136 (1.7%)	93959 (1.7%)	1,177 (1.6%)
<b>Smoking Status</b>			
Non smoker	5230971 (93%)	5166174 (93%)	64697 (86%)
Current or former smoker	391431 (7.0%)	381036 (6.9%)	10395 (14%)
<b>Sex</b>			
Females	3231360 (57%)	3195925 (58%)	35435 (47%)
Males	2390942 (43%)	2351285 (42%)	39657 (53%)
<b>Race/Ethnicity</b>			
Hispanic or Latino Any Race	688110 (12%)	679102 (12%)	9008 (12%)
Black or African American	661217 (12%)	646374 (12%)	14843 (20%)
Non-Hispanic			
White Non-Hispanic	3403103 (61%)	3359522 (61%)	43581 (58%)
Other Non-Hispanic	869872 (15%)	862212 (16%)	7660 (10%)
Comorbidities > 4	152187 (2.7%)	147787 (2.7%)	4400 (5.9%)



# Executive Summary

## Cluster 1:

### Low-risk

Lowest age  
Hispanic or Latino  
Other Non-Hispanic

No or lowest number  
of comorbidities

Relative healthier  
indicated by lab-  
related variables

## Cluster 2:

### Moderate-risk

Female

Older age

White Non-Hispanic  
Lowest BMI

Highest CD4 count  
percent

## Cluster 3:

### High-risk

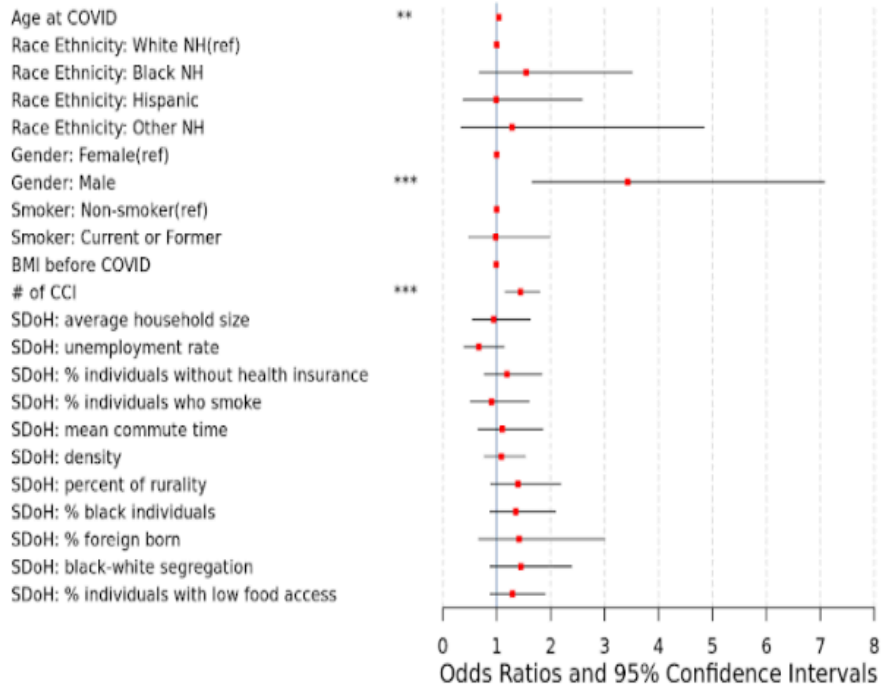
Male  
Current or former smoker  
Black Non-Hispanic  
Highest BMI  
Highest HIV viral load

Most number of  
comorbidities  
Higher proportions of death,  
long COVID &  
hospitalization

Cluster Comparison	Odds of COVID-19 Severity	Odds of Death	Odds of Long COVID-19
Cluster 3 vs Cluster 1	20 times more likely	11 times more likely	1.3 times more likely
Cluster 2 vs Cluster 1	3 times more likely	2 times more likely	Similar odds
Cluster 3 vs Cluster 2	6 times more likely	6 times more likely	1.23 times more likely

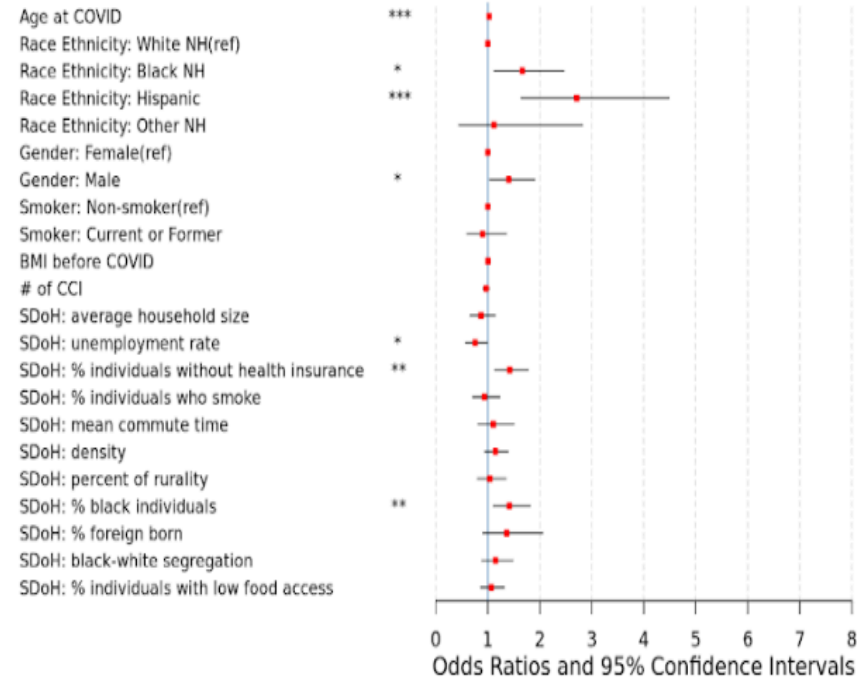
# Cluster 1 vs Cluster 3- Severe COVID-19

## Cluster 1: Low-risk cluster



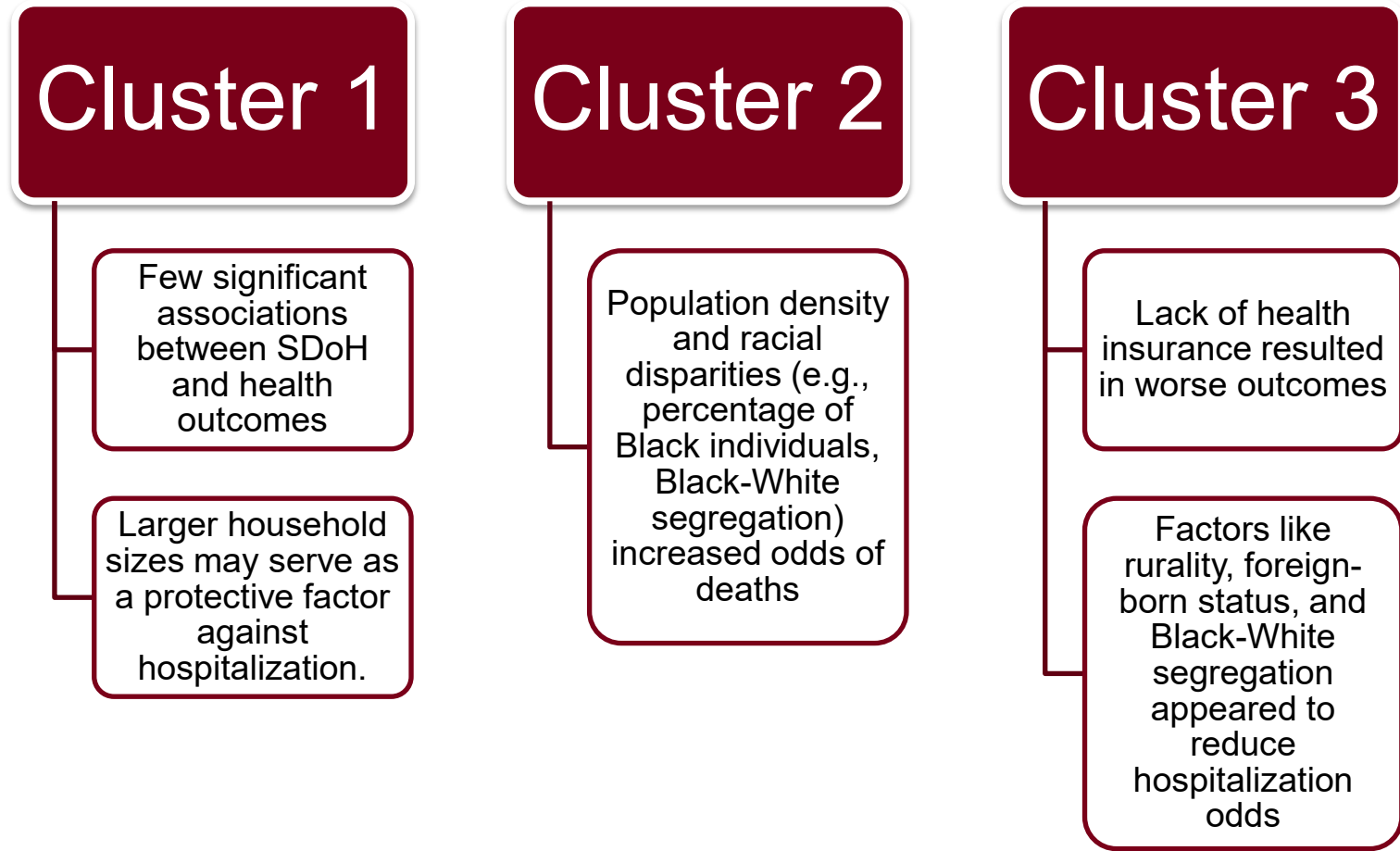
← Less likely to develop severe COVID-19      More likely to develop severe COVID-19 →

## Cluster 3: High-risk cluster

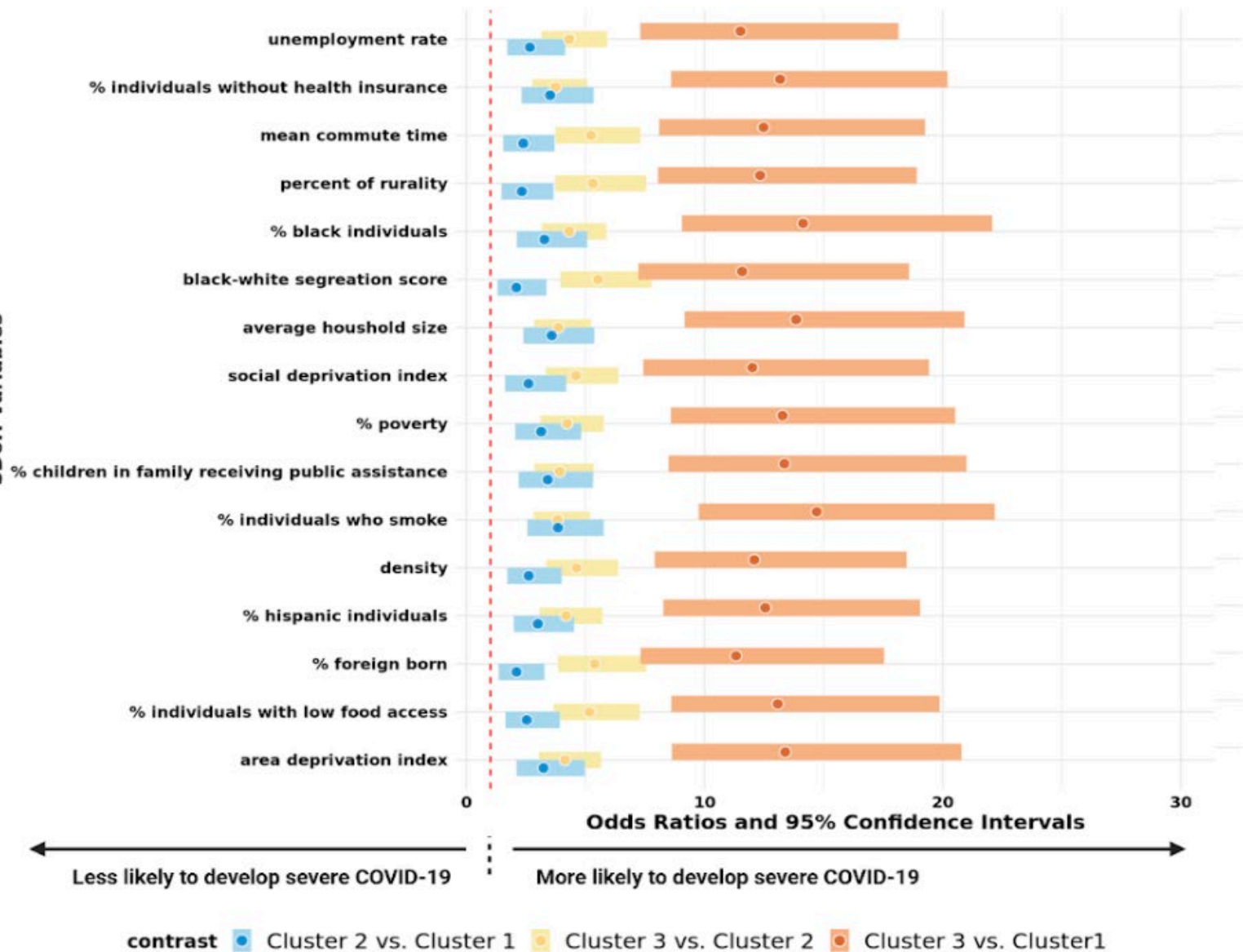


← Less likely to develop severe COVID-19      More likely to develop severe COVID-19 →

# Comparison of SDoH across Clusters



# SDoH Variables



# Implications of findings



Offer deeper insights into characteristics of PLWH with worse COVID-19 outcomes



Could enable targeted policy interventions and improve health equity



Provides a deeper perspective of the intersections of social determinants of health and COVID-19 clusters in PLWH



Approach could be adopted to other clinical outcomes in HIV



# Limitations of Work

---

Missing data in lab measurements

---

Resampling approach might affect our findings

---

Limited sample size for HIV-specific variables

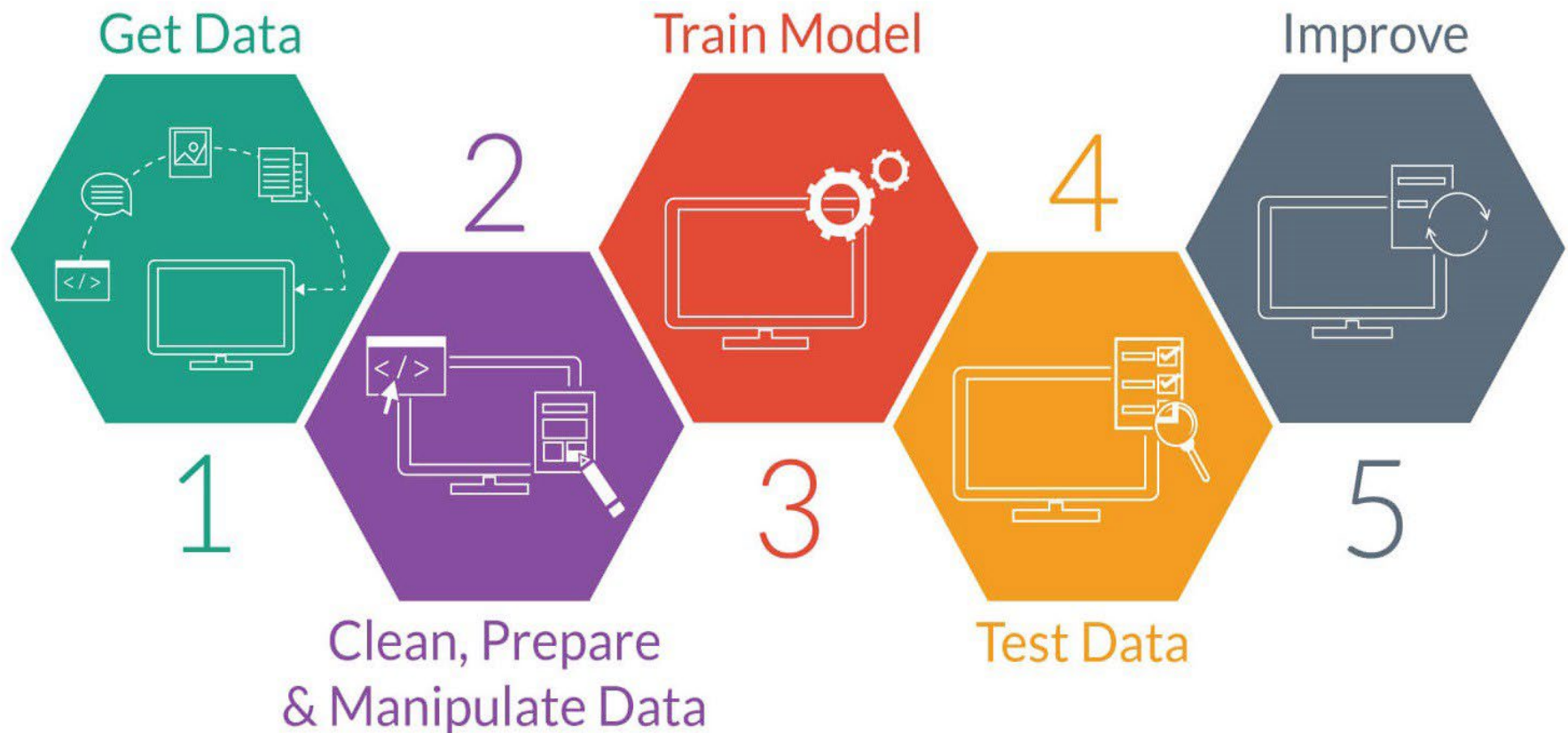
---

No adjustment for vaccination





# How can you get started with ML?



[This Photo](#) by Unknown Author is licensed under [CC BY-SA-NC](#)

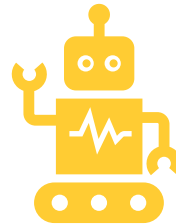
# Some Challenges and Considerations of ML



## Ethical concerns

Data privacy (HIPAA-compliance)

Algorithmic fairness



## Limitations of ML

Clinical judgment vs machine learning predictions

The need for transparent, interpretable models (e.g., why a prediction was made)



# Conclusions and Takeaways



## Summary

ML is an exciting tool that can improve healthcare, particularly in EHR

Can improve precision medicine, detect disease early, reduce human error, identify patients for clinical trials and many more



## Call to action

Consider exploring simple ML models in your work

Get involved with others using ML and learn to do the dirty work

Learn to code in R/Python (Coursera etc.)

I just need  
the main ideas



**Thank You**

---

Tiankai Xie

---

N3C HIV-Subdomain team

---

NIH Funding







UNIVERSITY OF MINNESOTA

**Driven to Discover®**

Crookston Duluth Morris Rochester Twin Cities

The University of Minnesota is an equal opportunity educator and employer.