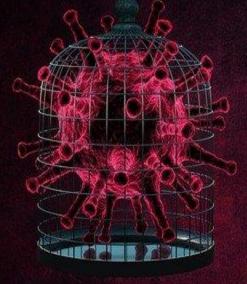
# **Covid-19 Severe Outcome Risk Prediction**



Private Machine Learning on Medical Records & Social Data

Changrong Ji

Dr. Mahesh Shukla

Dr. David Patton

Dr. Xue Yang

Dr. Xingguo Zhang

Antonio Linari

**Premdutt Gaur** 

Vance Degen

**A3.AI** 

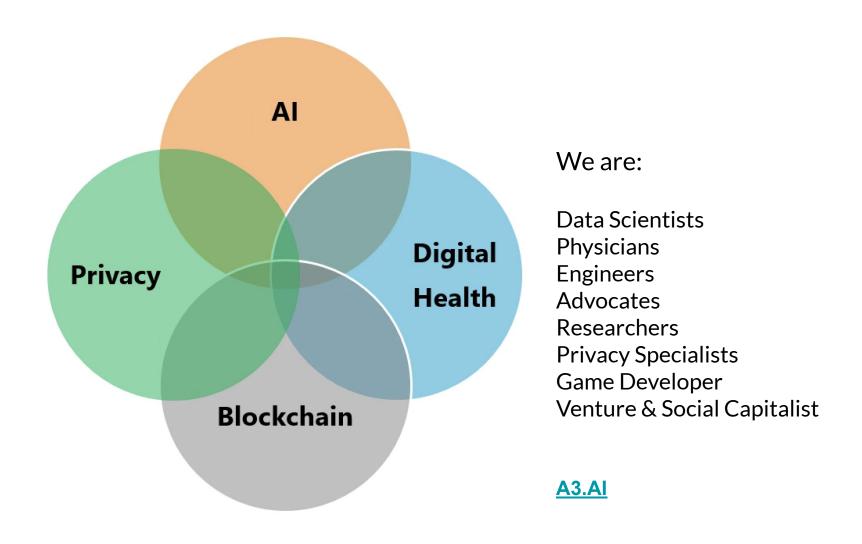
# **Topics**

- About Us
- Aims
- Data
- Approach
- Early Findings
- Future Work





# Nonprofit Applied R&D





# **Projects**



Covid-19 Severe Outcome Risk Prediction

Privacy-preserving Machine Learning on Medical Records & Social Data

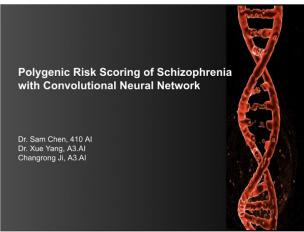
Changrong Ji
Dr. Mahesh Shukla
Dr. David Patton
Dr. Xue Yang
Dr. Xinguo Zhang
Antonio Linari
Vance Degen
A3.AI



Private & Secure Al



**Health Analytics** 



**Drug Safety & Repurposing** 



Reinforcement Learning

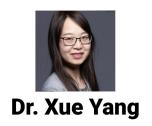


#### **COVID-19 Project Team**



















#### Aims

1. COVID-19 Severe Outcome Risk Prediction

2. Social Determinants of Health and Risk Factors of COVID-19

3. Privacy-preserving Machine Learning

 Building Clinical Concept Embeddings under Computing Resource Constraints



Multiple research aims are addressed in detail in the following working papers respectively as of 09/2020. Future versions will be published as the research progresses:

- 1. Early findings from machine learning baseline models to predict an individual's risk of hospitalization if infected with COVID-19, based on medical claims and EHR, respectively: <a href="Predicting Risk of Hospitalization among COVID-19 Patients">Predicting Risk of Hospitalization among COVID-19 Patients</a>, Mahesh Shukla et al.
- 2. <u>Social Determinants of Health for COVID-19 Diagnosed Patients</u>, David Patton et al.
- 3. The first in a series of private AI techniques which will enable the release of the AI models built with personal level data while preserving privacy: <a href="Privacy-Preserving Machine Learning Techniques 2020">Privacy-Preserving Machine Learning Techniques 2020</a>, Changrong Ji et al.
- 4. Clinical Concept embedding is a feature engineering technique to enhance the accuracy of AI models. To build embeddings from large claims data typically requires high computing power. We use a novel approach to efficiently build embeddings under resource constraints in the COVID-19 Research DB environment: <a href="Building Clinical Concept Embeddings under Computing Resource Constraints">Building Clinical Concept Embeddings under Computing Resource Constraints</a>, Antonio Liniari et al.

#### Data

As of 08/21/2020, with new data added with 1 week delay

#### Claims

- 98 million patients 7 years of medical claims history of over 3 billion claim lines
  - Key attributes: ICD diagnosis codes, CPT procedure codes
- o 200,000+ COVID patients

#### • Electronic Health Record (Outpatient)

- 36 million patient's outpatient EHR records;
  - Diagnoses, Procedures; Encounters, Medications; Allergy, Social History, etc.
- o 16,000 confirmed COVID patients, and 75,000 possible COVID patients

#### Social - Claims - Death linked

- 242 million people's Social Data
  - People (demographics, finance, credit, housing, jobs, lifestyle)
  - Behaviors (interests, purchasing, social network activity, charitable giving, health lifestyle)
  - Predictors (motivator, travel, auto, in-market, and economic stats)
- Death Registry of 80% of US population
  - Died in 2020
- 95,000 COVID patients



#### Attributions to Data Providers

#### AnalyticsIQ

AnalyticsIQ is s a leading predictive data and analytics innovator that leverages a blend of publicly available data and custom algorithms informed by cognitive psychology concepts to describe consumers across three areas - People, Behaviors, and Predictors. Headquartered in Atlanta and recently named one of Georgia's Top 10 most innovative companies, AnalyticsIQ's team of data analysts, scientists, and cognitive psychologists have over 100 years of collective analytical experience and expertise.

#### Health Jump

Electronic Health Record data including diagnosis, procedures, labs, vitals, medications and histories sourced from participating members of the Healthjump network.

De-identified claims data was contributed by a claims clearinghouse.

# **Machine Learning for Clinical Prognosis**



#### Aim 1

Create machine learning models to predict a patient's risk of severe clinical outcomes if infected with COVID-19.

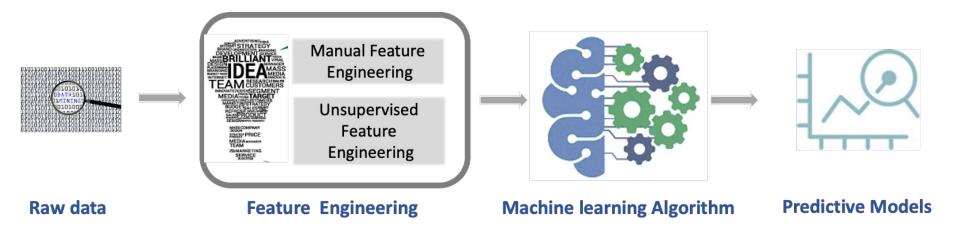
- Hospitalization
- ICU
- Intubation
- Ventilation
- ECMO (heart-lung bypass)
- Death, etc

These personalized risk scores and associated risk factors analysis can

- Help citizens make informed work and lifestyle choices
- Augment clinical prognosis by physicians
- Help health care organizations coordinate care and optimize resources
- Help public health agencies with planning, responding and reopening.

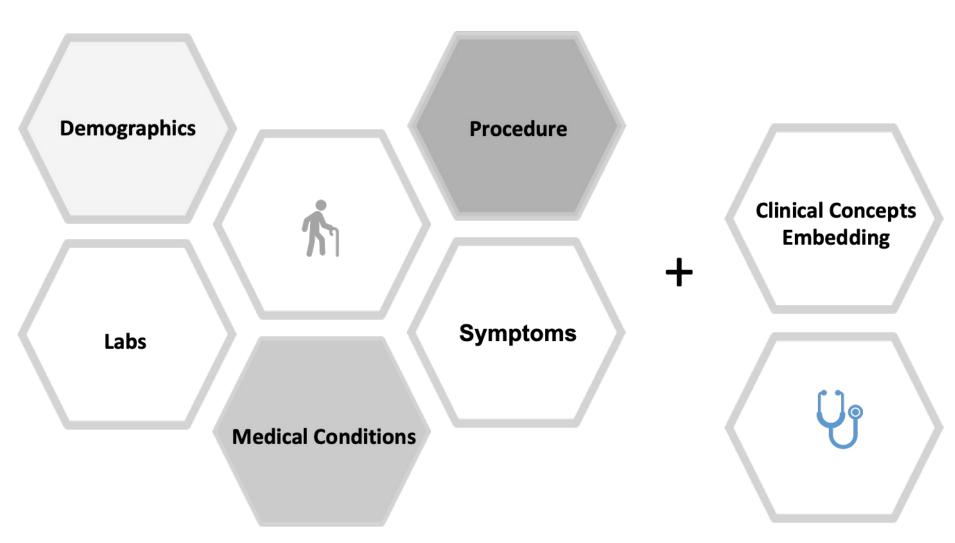


# ML Model Development





# Feature Engineering





# Embedding in NLP

A great way to represent sparse, high-dimensional data in NLP

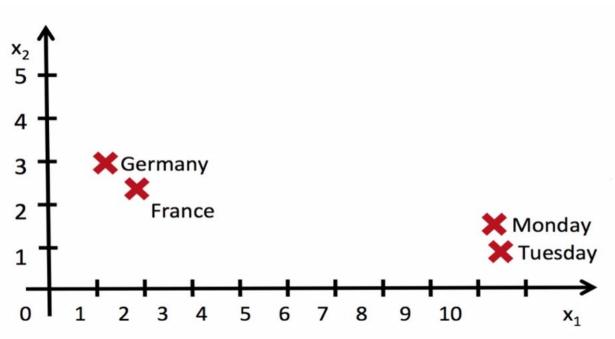


Figure (edited) from Bengio, "Representation Learning and Deep Learning", July, 2012, UCLA

#### Lower-dimensional space:

- Words of similar meaning are located near each other in the embedded vector
- Relative location of two words in the space could encode a meaningful relationship



# Clinical Concepts Embedding



Low dimensional continuous representation of high dimensional discrete data (~300,000 distinct codes).



One type of pretrained model used for transfer learning



# Clinical Concepts Embedding

Low-dimensional vector representations of medical concepts

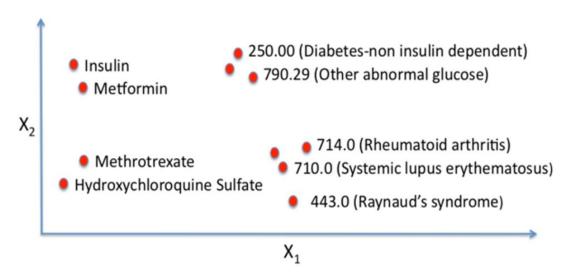


Figure 1: Illustration a low-dimensional representation (in this case, 2 dimensions) of medical concepts. Similar concepts are close to each other in Euclidean space.



# Hospitalization 15.7% ICU 6.5% Ventilation 0.4% ECMO (heart-lung bypass) 0.01% Death, etc 3.1%

# Data Labeling



# Top 20 Procedures for COVID Patients

CLM + CLM_INST+SVC+SVC_INST	COVID cases 98,182	7/19/2020		
PROCEDURE_CODE	Number of patients	Procedure description		
U0003	13915	COVID testing using high-throughput technology		
99213	11707	1707 Established patient office visit 0967 ER visit for E&M 9432 SARS Cov-2 Antibody test		
99285 86769 99233 93010 99223 99291 99232 71045	10967			
	9432			
	8666	6 Initial hospital inpatient care		
	7938	ECG		
	7498	NULL		
	7390	Initial hospital inpatient care		
	7332	Critical care, first hour		
	7254	Subsequent hospital care		
	6456	Diagnostic imaging chest		
99284	4900	ER visit for E&M		
87635	4463	SARS Cov-2 DNA RNA test		
99283	4436	ER visit for E&M		
99309 99212 85029 99214	4393	Subsequent nursing facility care		
	4271	Office visit E&M		
	4098	Lab automated diff wbc count		
	3619	Office visit E&M		
80053	3296	Compr metabolic panel		
36415	3243	Collection of venous blood		
99308	2737	Subsequent nursing facility care		



# Top 20 Co-occurring Diagnosis with COVID-19

A		t .	
Diagnosis code	#patients	Diagnosis description	
J1289	12321	Other viral pneumonia	
110	10671	Essential hypertension	
J9601	8078	ARF with hypoxia	
J189	7885	Pneumonia unspecified organism	
RO5	7118	Cough	
R0602	6165	Shortness of breath	
E119	5375	Diabetes Mellitus	
R509	4749	Fever, unspec	
R0902	4325	Hypoxemia	
N179	4073	AKF, unspeci	
E785	3271	Hyperlipidemia, unsp	
A419	3185	Sepsis, unspecified org	
Z20828	2832	Contact	
N390	2652	UTI, site not sp	
M6281	2283	Muscle weakness gene	
J988	2235	Other sp resp disorders	
D649	2228	Anemia, unsp	
R531	2063	weakness	



#### Baseline Prediction with Claims Data

**Goal**: Prediction of hospitalization for COVID-19 patient

#### Data:

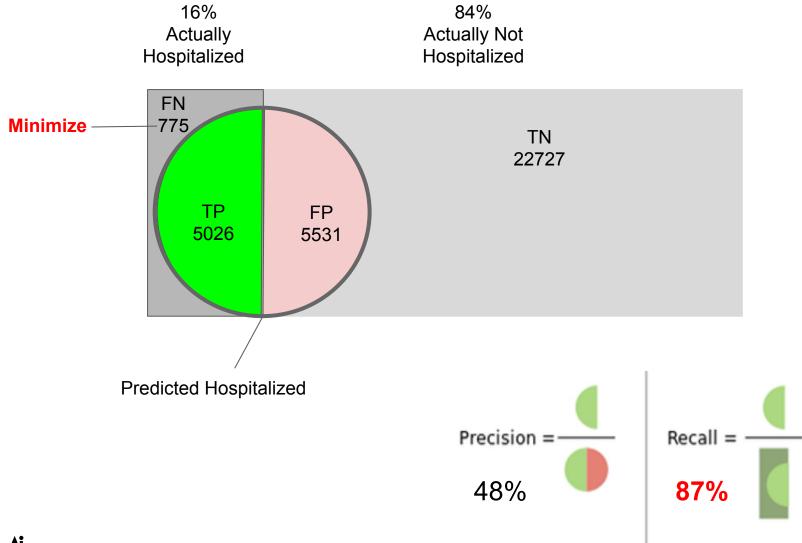
- Hand-crafted features (~100)
- Total patients with COVID-19 infection: 170,241
- Hospitalized patients: 17% of above
- Train/Test/Validation set split: 60/20/20

#### Model:

- Random Forest Classifier
- Balanced weights
- 400 estimators with a max\_depth of 20



## Precision & Recall Refresher





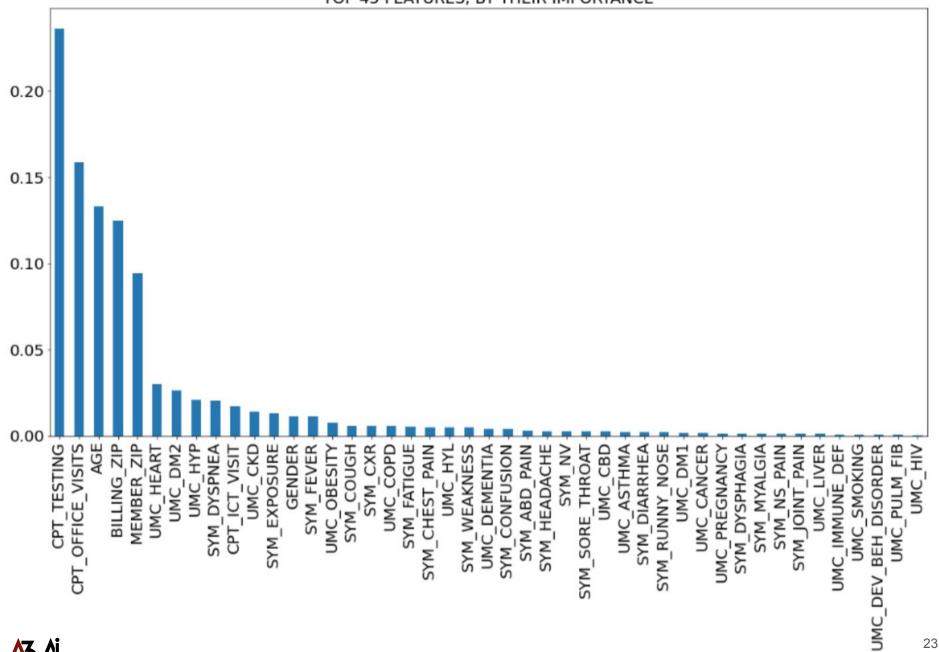
# Hospitalization Results

## Classification Report:

		precision	recall	f1-score	support
	0	0.97	0.80	0.88	28248
	1	0.48	0.87	0.61	5801
accuracy				0.81	34049
macro av	g	0.72	0.84	0.75	34049
weighted av	g	0.88	0.81	0.83	34049



#### TOP 45 FEATURES, BY THEIR IMPORTANCE



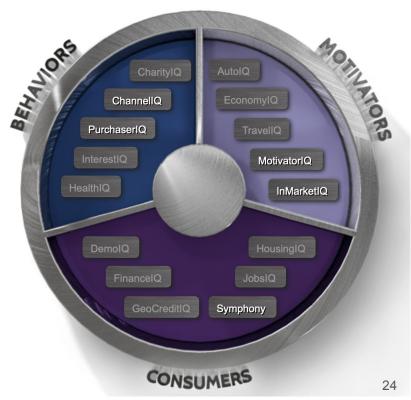


#### Social Determinants and Risk Factors

About 90 attributes of social data from Analytics IQ are available for over 34 million patients. Over 95,000 are COVID-19 patients.

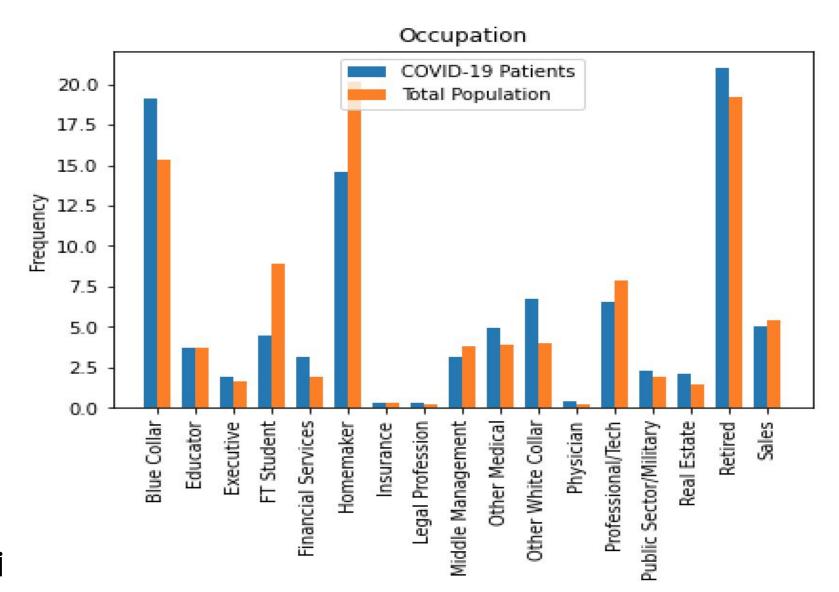
#### We examined:

- Impact of different demographics
- Behaviors as risk factors
- Predictors (likelihood) as risk factors



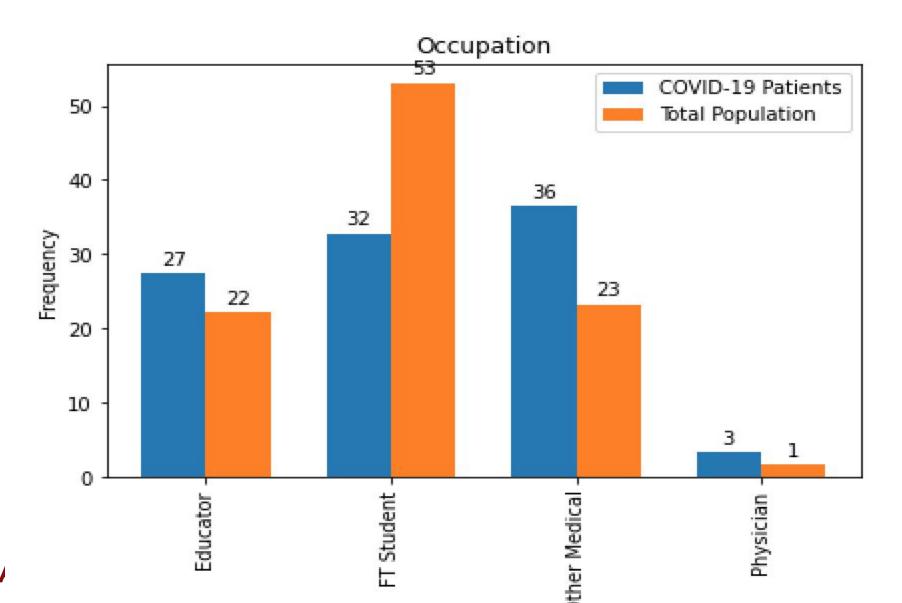


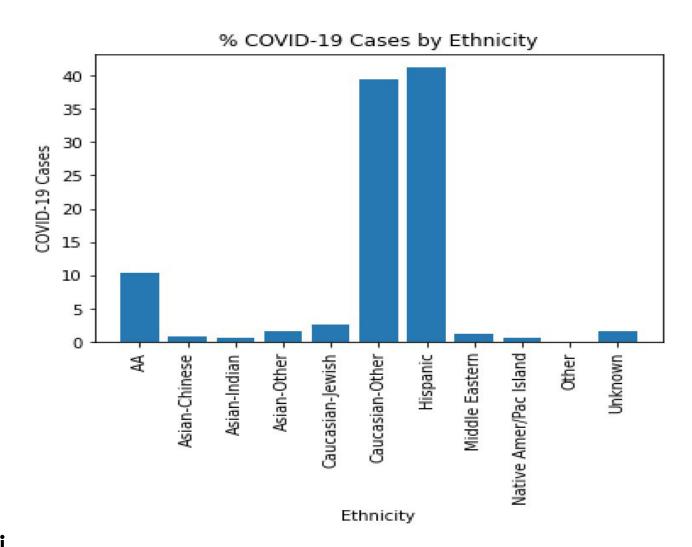
## Occupation\*



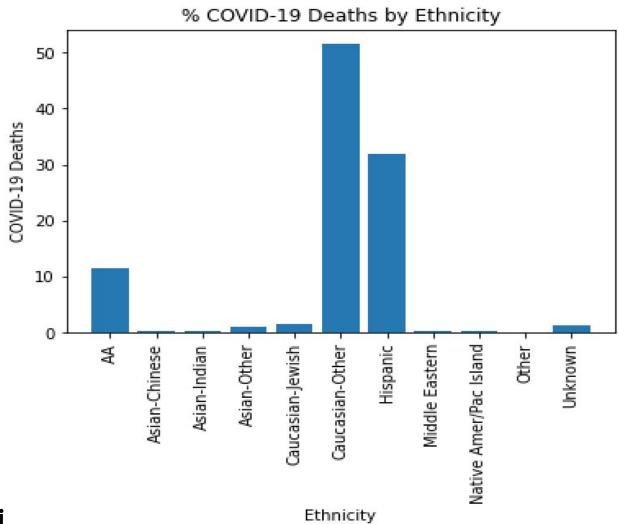


# **Specific Occupations**

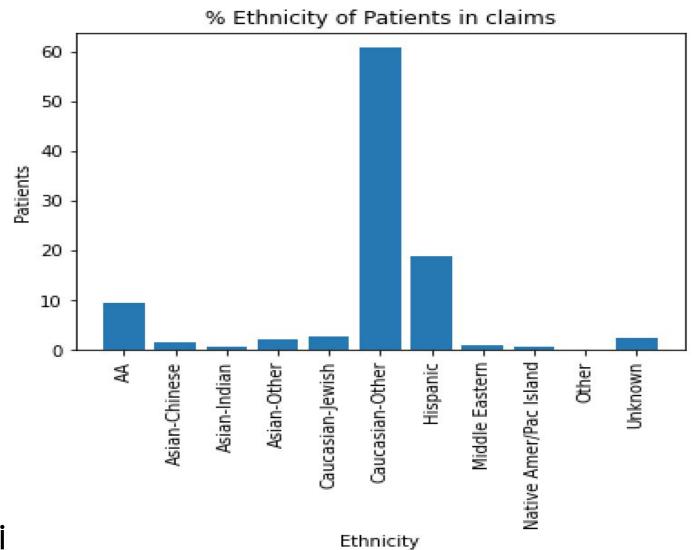




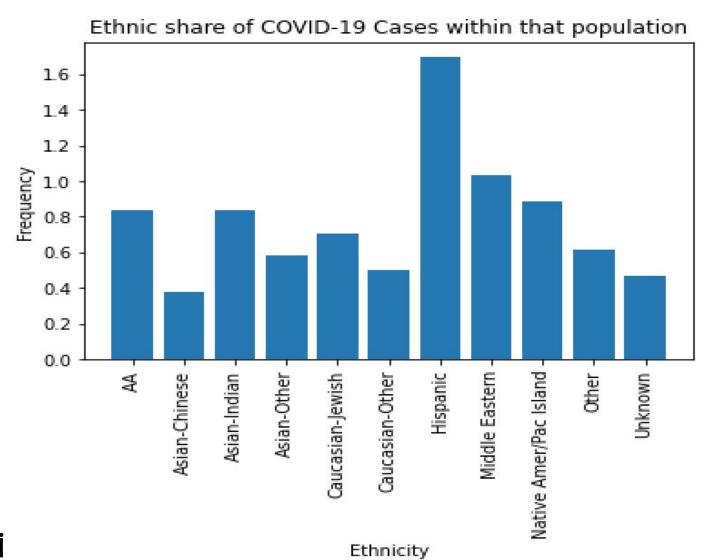






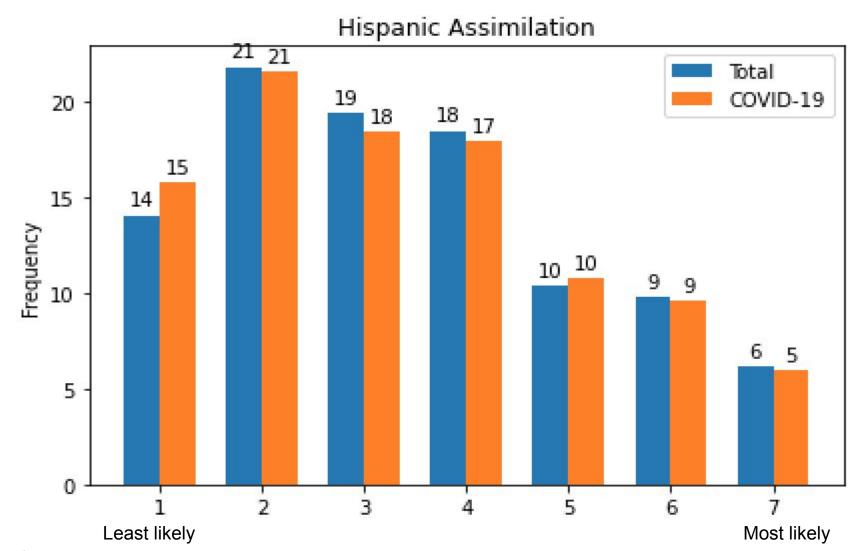








## Assimilation into US Culture





## **Future Work: SDOH**

Many additional attributes available in the dataset:

BMI, Diet, Location, Profession, Access to Healthcare, behavior, etc

- Population Health study
- Predictive Models
- Personalized risk scoring



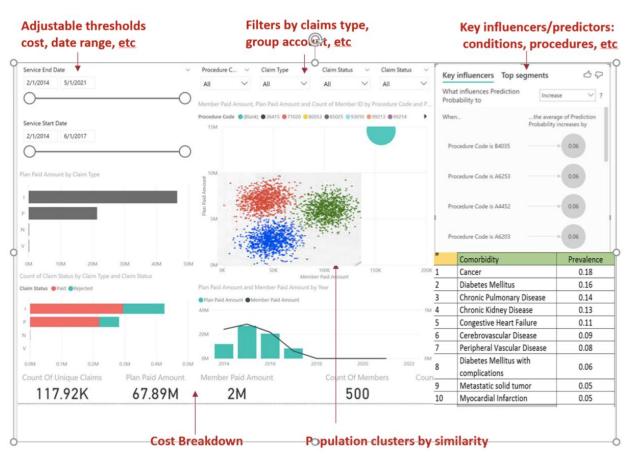


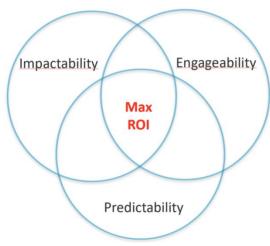
#### Future Work: Current Use Case

- Improve model performance
  - Add clinical concept embeddings
  - Add temporal features (time sequence of clinical events)
  - Other machine learning and (lighter weight) deep learning models
- Add more classification categories (other severe outcome types)
- Incorporate Electronic Health Record data for prediction



## Future Work: Population Health Dashboard







## Future Work: New Use Cases As the Pandemic Progresses

- Drug and COVID-19 vaccine effectiveness and safety
- Identify clinical trial candidates
- Long term consequences of COVID-19
  - Personal
  - Societal
  - Health
  - Economical





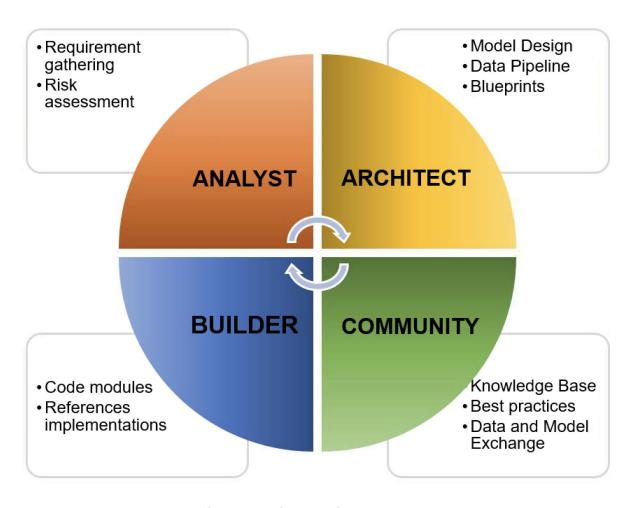


## Data Sharing & Healthcare Al Challenges

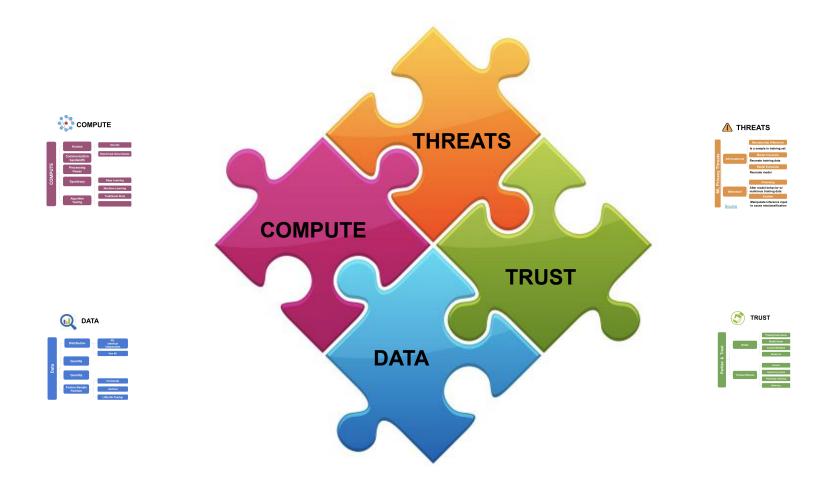
- Advanced analytics relies on data, often curated from multiple sources.
- Data sharing and usage have a complex web of trust with multiple parties: data owners, analytics solutions providers and users.
- Aggregating data and building models centrally raise concerns in
  - privacy and security,
  - single point of failure,
  - intellectual property and
  - o misaligned incentives on the usage and value of of the combined data and models
- CLOAK is a platform and toolkit that is under development to address some key challenges in collaborative privacy preserving machine learning.

Specific relevance to the COVID-19 Research DB projects: The personal level medical records and social data, while de-identified, are still vulnerable to attacks such as data linkage and model inversion that leaks private information. The following highlights a set of techniques to mitigate the privacy risks. A series of papers will be published on this topic. Starting with: <a href="Privacy-Preserving Machine">Privacy-Preserving Machine</a> <a href="#arning Techniques 2020">arning Techniques 2020</a>, Changrong Ji et al.

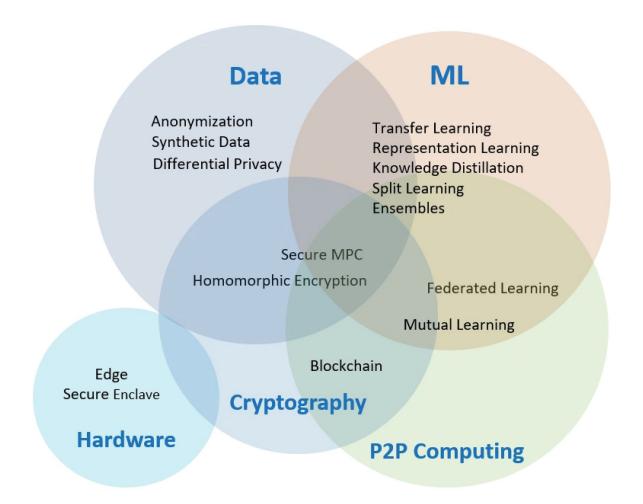
# **CLOAK PLATFORM (work in progress)**



## **ANALYST**



## **PRIVACY PRESERVING TECHNIQUES 2020**

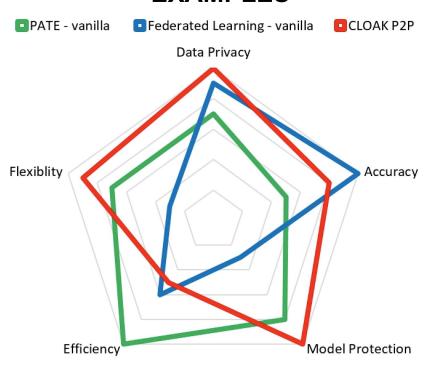


## **ARCHITECT**

#### **DESIGN TRADE-OFFS**

# Solution Flexibility Computing Efficiency Model Protection

#### **EXAMPLES**



Copyright © 2020 Changrong Ji

## **BUILDER**

