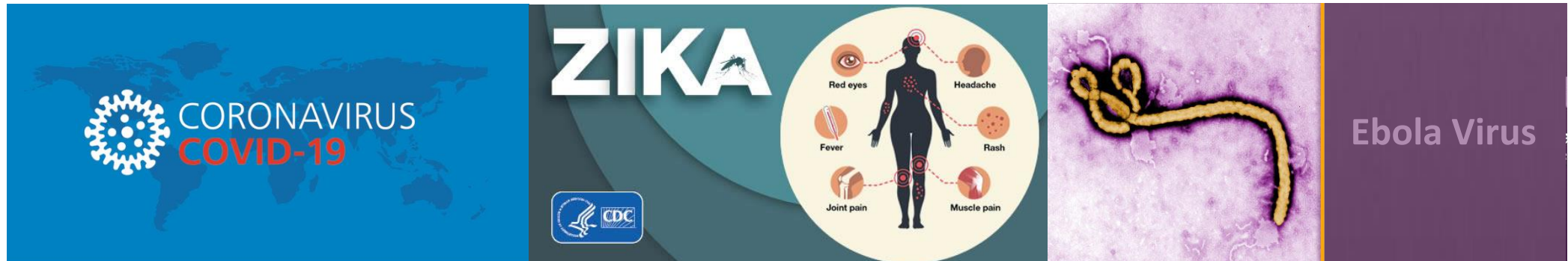# Global Health Challenge

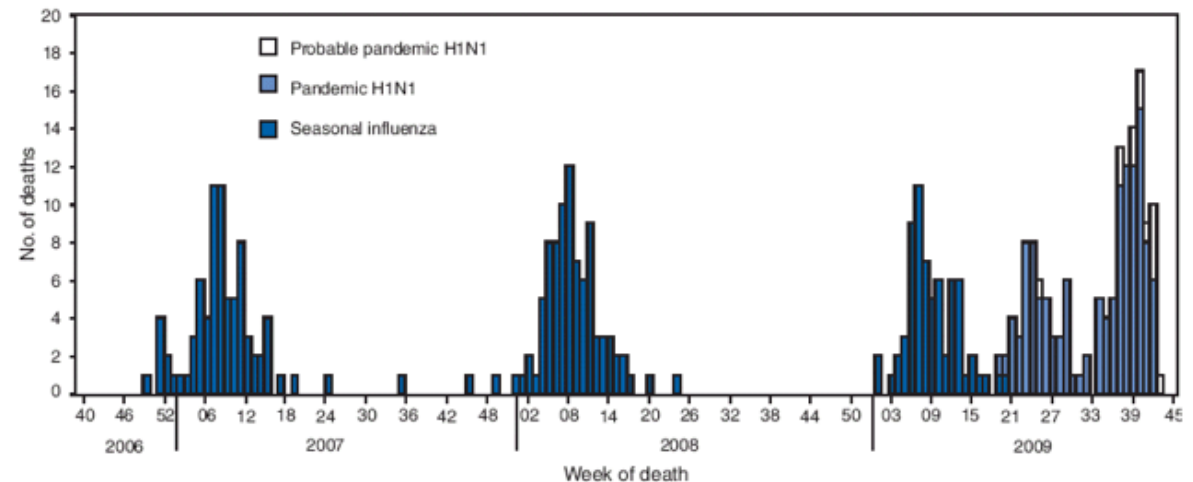DATA SCIENCE IN INFECTIOUS DISEASES

# Background

- Infectious diseases are closely related to our human world.

- Influence everybody's life in a certain way.

- Travel across continents and spread widely.



CORONAVIRUS COVID-19

ZIKA

Red eyes    Headache
Fever    Rash
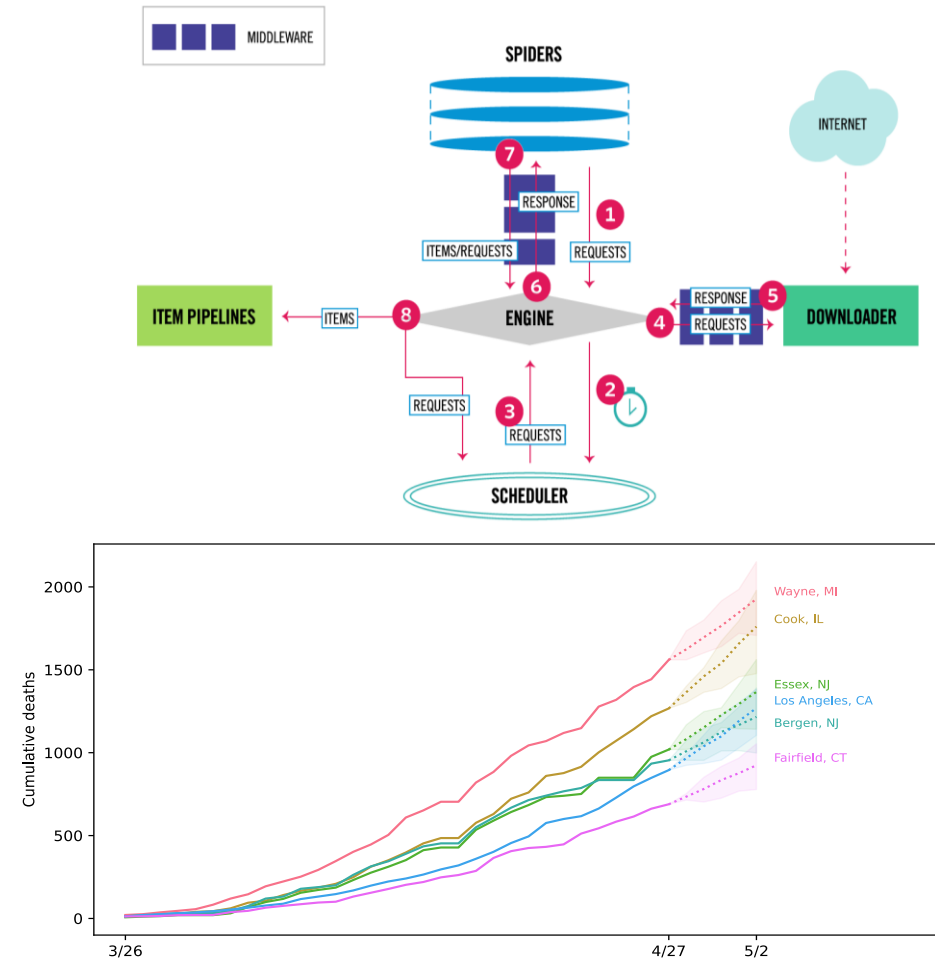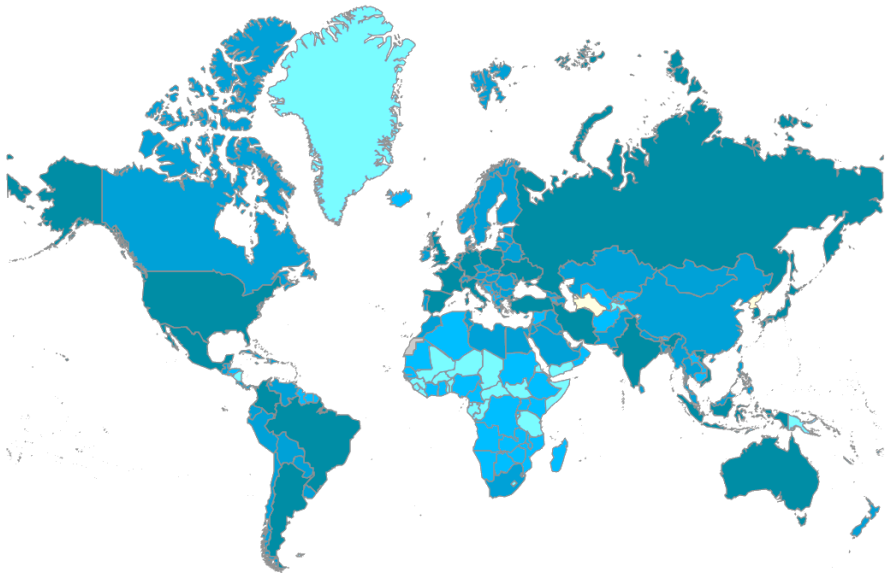Joint pain    Muscle pain

CDC

Ebola Virus

# Technical Ideas

- Infectious diseases have extensive existences.

- Carefully monitored and studied by scientists around the globe.

- Many of their pandemics take place periodically.

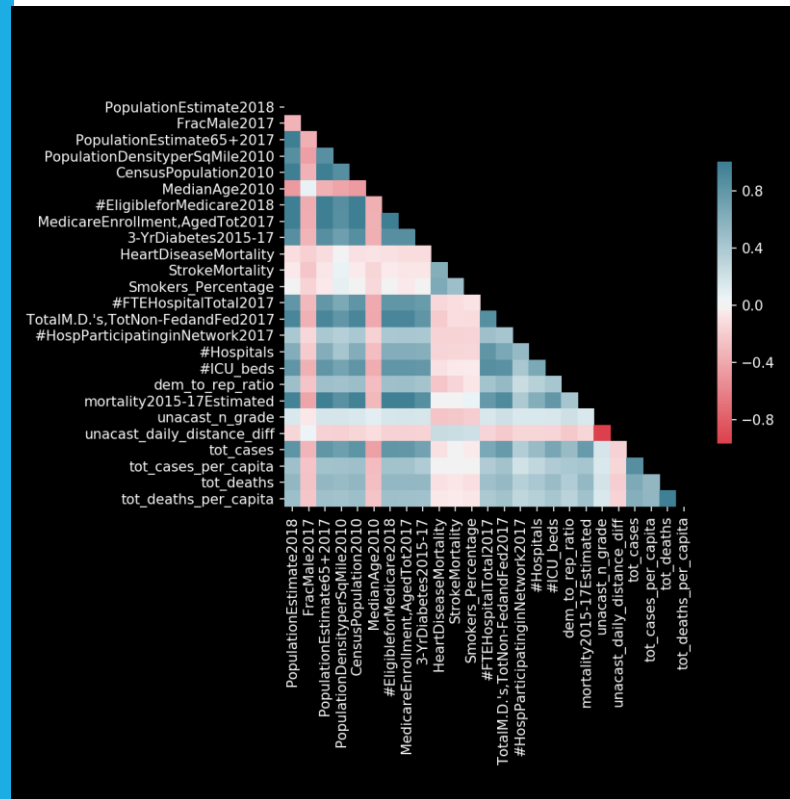- Many of their spread and mutation are strongly featured.

# Technical Route

- Sampling and collecting data from the Internet.

- Visualization using the collected data.

- Data mining and further analysis.

# Sampling and Collecting Data

- [ ] High dimensional data from multiple sources on the Internet.

- [ ] Automatically update and parse data using web spiders.



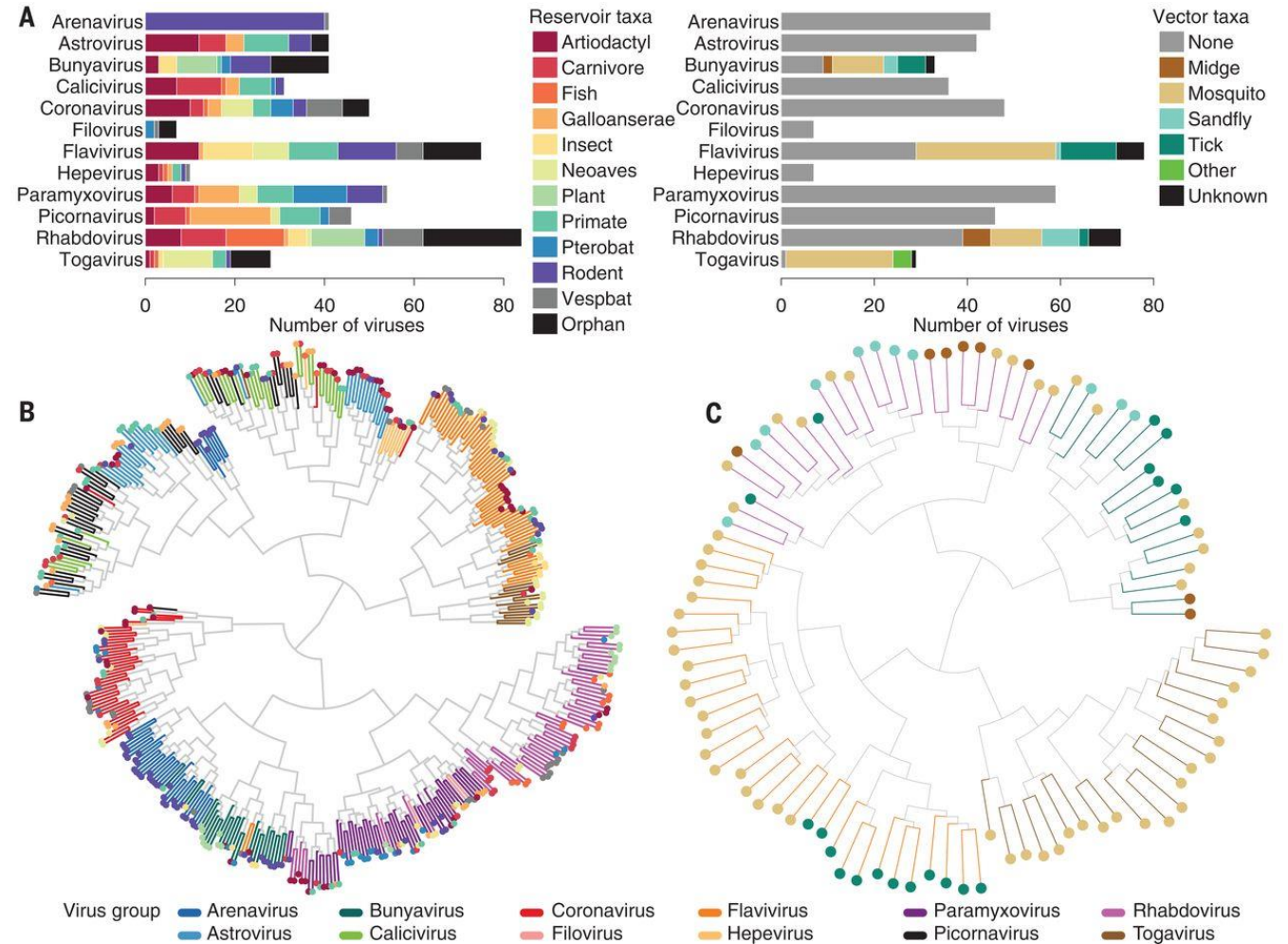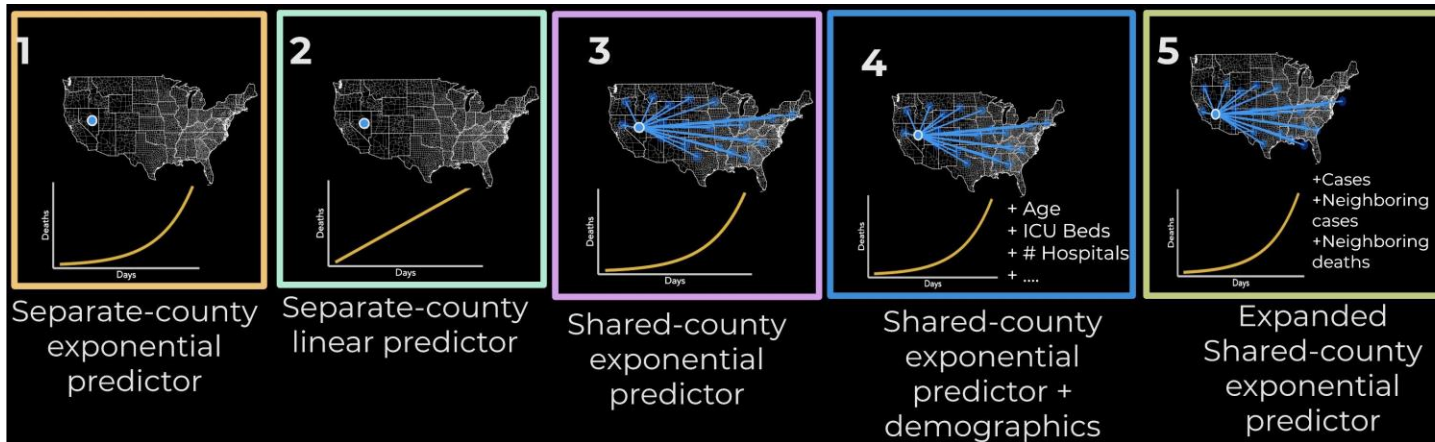| Month | Actual value | Predicted value | Absolute error | Percent absolute error |
|-------|-------------|-----------------|----------------|------------------------|
| 01–2012 | 10046 | 10230 | 184 | 1.8% |
| 02–2012 | 17421 | 14578 | 2843 | 16.3% |
| 03–2012 | 21625 | 18429 | 3196 | 14.8% |
| 04–2012 | 10707 | 11785 | 1078 | 10.1% |
| 05–2012 | 8520 | 8618 | 98 | 1.2% |
| 06–2012 | 6195 | 6621 | 426 | 6.9% |
| 07–2012 | 6738 | 5240 | 1498 | 22.2% |
| 08–2012 | 6793 | 5983 | 810 | 11.9% |

# Data Visualization

- ☐ Generate charts from the data automatically.

- ☐ Provide clear and overall images of pandemics.

- ☐ Help scientists' research and public education.

# Data Mining and Analysis

- ☐ Learn the trend of infectious disease with epidemic models.

- ☐ Mine the data using Deep Learning Networks.

- ☐ Provide statistic data for other scientific researches.
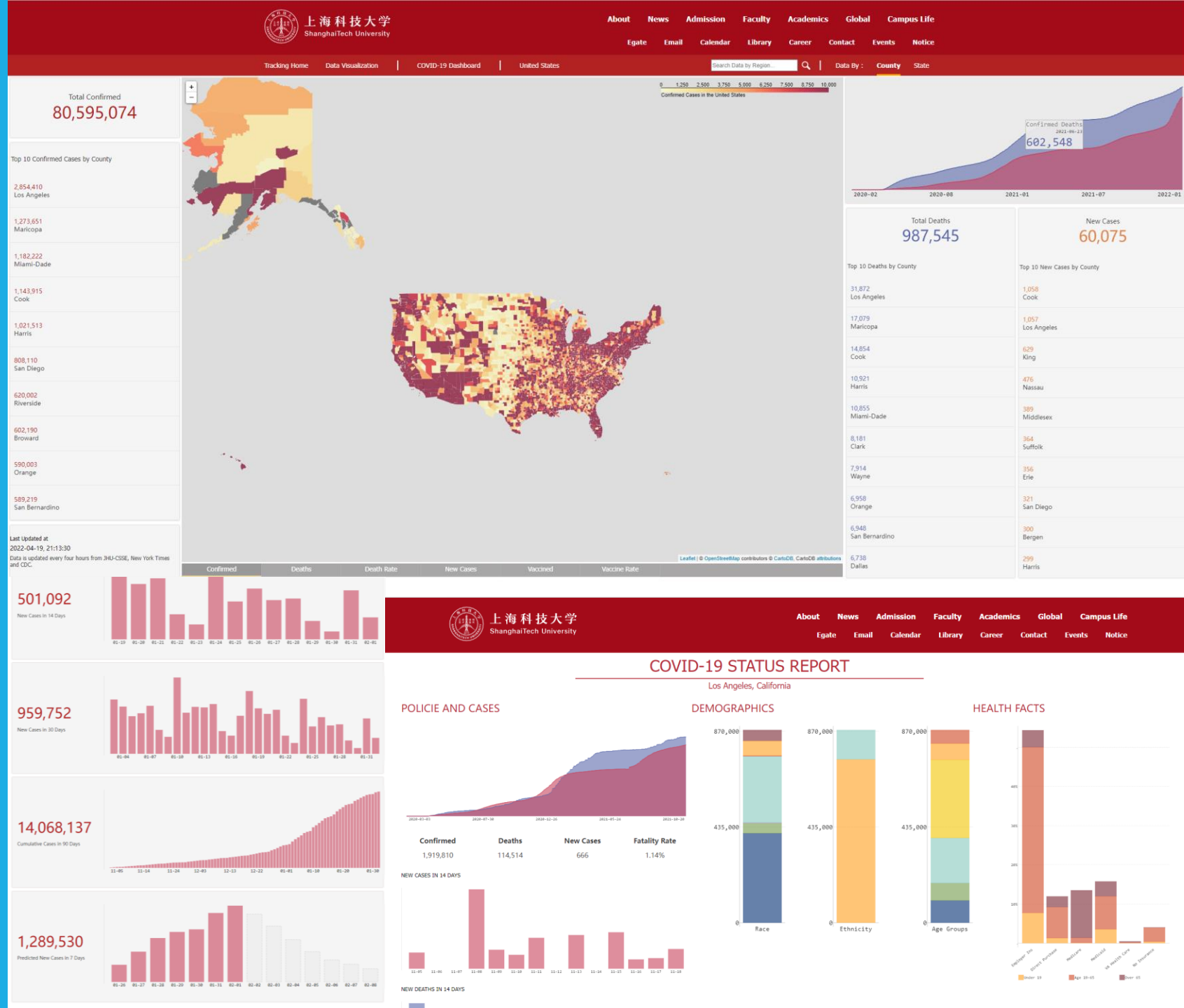
# Difficulties



- Inconsistent data form.

- Incomplete or broken data pieces.

- Useless or disturbing data.

- Customize Deep Learning Models.

- Select data with different features.

# Example

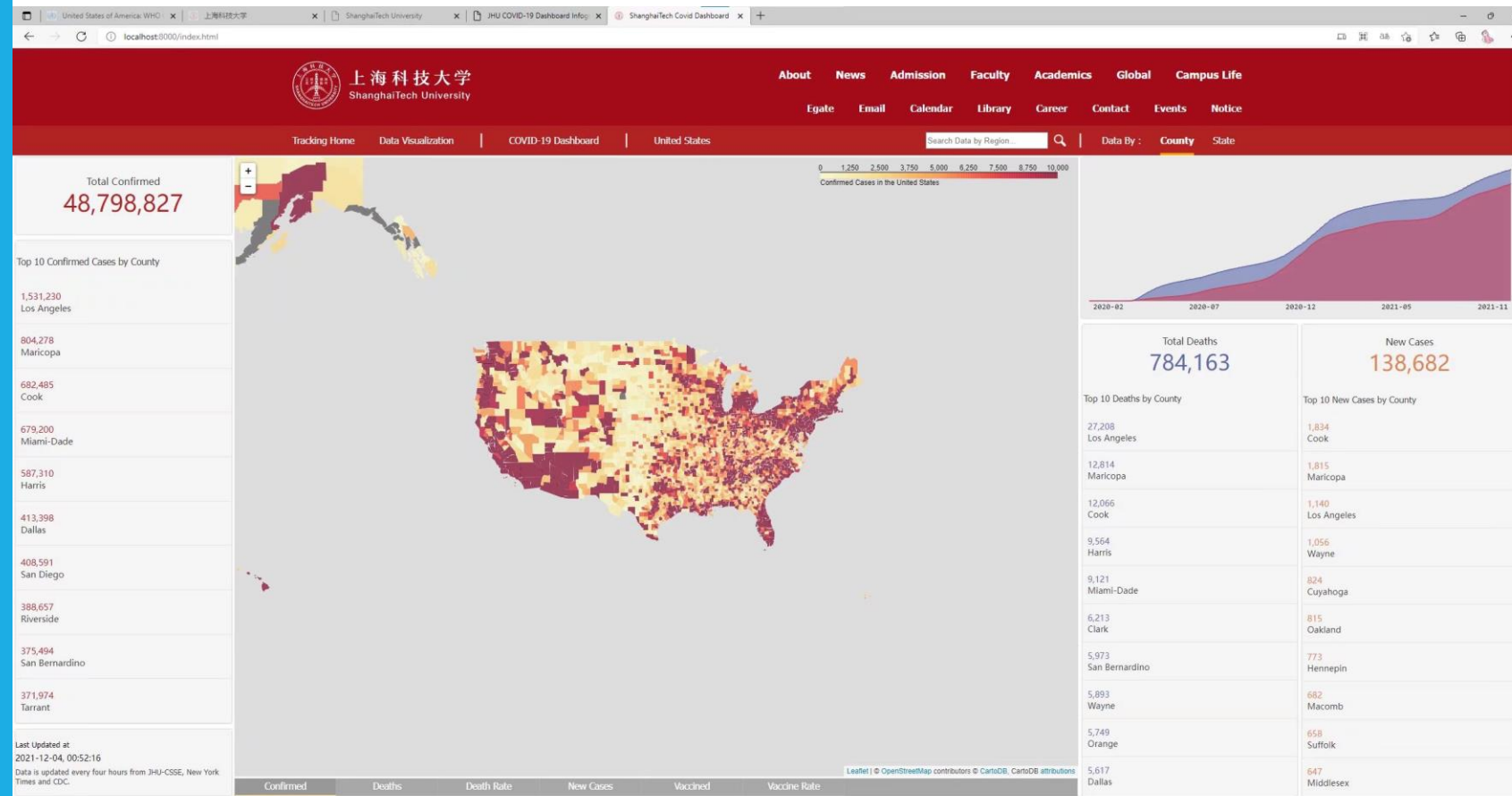A prototype build by our team.

Website ( in ShanghaiTech campu )
http://10.19.75.90:12345/index.html

Source code:
https://github.com/yanglinshu/covid
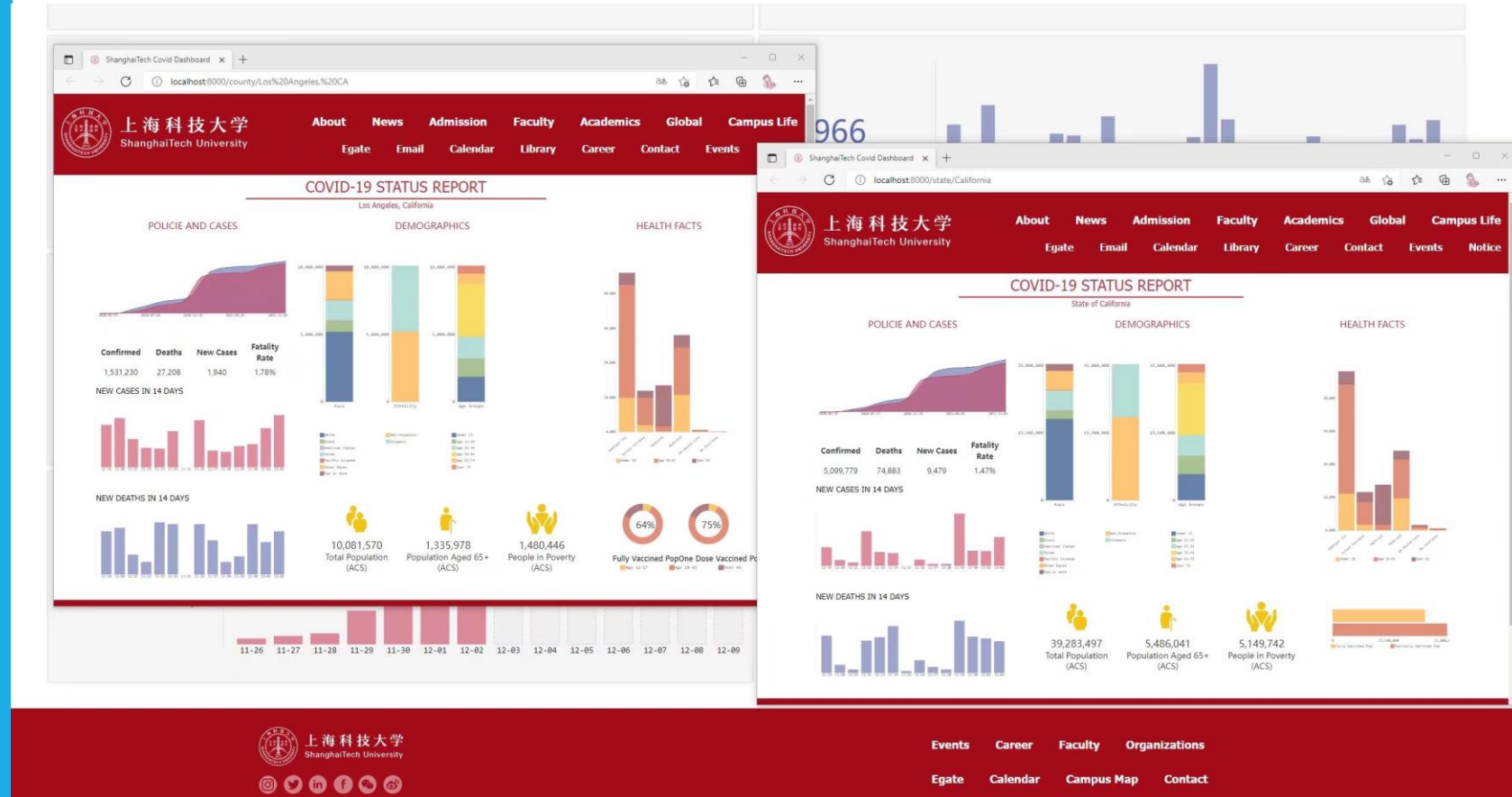
# Example

□ Our Data

collected from Johns Hopkins University, CDC, Census Bureau, New York Times, using a web spider based on python. Parsed and cleaned using Dataframe and Pandas.
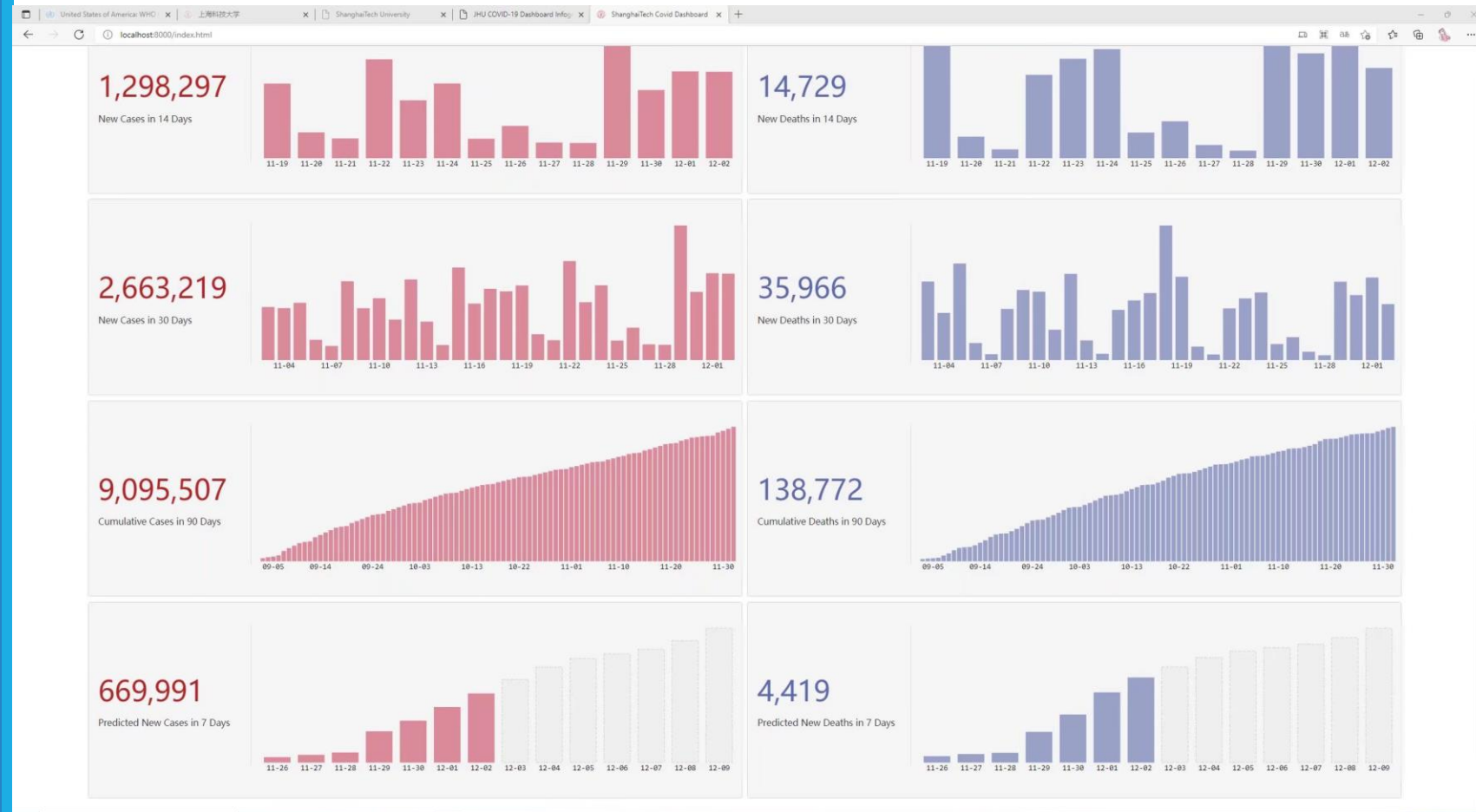
# Example

☐ Visualization

Generated using Pygal and Folium with geographical data from opendatasoft.com.

# Example

☐ Data Mining

A simple moving average model using the data of the last 14 days.

# Example

- [JieYingWu/COVID-19_US_County-level_Summaries: Attempt to find correlation between a region's demographic/economic factors with its ability to manage disease spread (github.com)](github.com)

- [Yu-Group/covid19-severity-prediction: Extensive and accessible COVID-19 data + forecasting for counties and hospitals. ☑ (github.com)](github.com)

- [facebookresearch/CovidPrognosis: COVID deterioration prediction based on chest X-ray radiographs via MoCo-trained image representations (github.com)](github.com)