# HFSSG: A Hybrid Framework combining SIR pandemic model and Spatio-temporal Graph neural network forecasting the COVID-19 Impact in the United States

Liqin Zhang[1], Siyin Ma[2], Zhejian Jin[3]

**Abstract**    Researches on infectious diseases has always been an important topic in interdisciplinary fields. The spread of infectious diseases not only depends on the attributes of diseases itself, but also depends on the network structure of the crowd. Since late 2019, COVID-19 has become a worldwide pandemic. In order to better understand and control the spread of COVID-19, many researchers collaborated and modeled this process to forecast its influence on the neighbouring areas. SIR model is a common compartmental model that well-defines disease transmission dynamics. However, it assumes a relatively closed system, and it does not take interactions among regions into consideration. Based on this, we utilize Graph Neural Network using data that describes inter-region interactions, which captures the features of messages among the regions. In this work, we propose a novel framework for COVID-19 case examination and prediction that uses Graph Neural Network. Different with existing time series models, this model learns from a spatio-temporal graph. We then evaluate this model on the California COVID-19 data set.

**Key words**    SIR(SEIR) Pandemic Model, Machine Learning, ST-Graph Neural Network

1. Student ID: 517370910123, Email: graves_zhang@sjtu.edu.cn
2. Student ID: 517370910003, Email: msy841@sjtu.edu.cn
3. Student ID: 517370910167, Email: jinzhejian@outlook.com

# Contents

# Introduction

## Problem Statement

With the rapid growth and spread of Covid-19, the ability to accurately forecast caseload is extremely important to help inform policymakers on how to provision limited healthcare resources, rapidly control outbreaks, and ensure the safety of the general public. In order to prepare, understand, and control the spread of the disease, we have come together in a collaborative effort to model and forecast COVID-19. In this project, we combine two powerful models (SIR and ST-GNN) into one integrate framework to perform better prediction capabilities in COVID-19 cases with lower cost. We propose a novel hybrid multi-network based framework named HFSSG that learns to select relevant edges and graph representations by pre-computing results from SIR model, along with RMSLE and Pearson values. On top of that, this model, with the help of ST-GNN, it learns the parameters for the pandemic model, and better fit the coronavirus model in real life. This main flowchart of our framework is shown in Figure 1.



**FIGURE 1.** *Flowchart of our main framework*

## Planned Approach

**SIR Pandemic Model**  SIR model is a simple method describing the epidemic spread that can be applied to any network structures. It divides the popularity into three categories, representing three stages of a person. The three functions with time $t$ is used to simulate the three compartments.

1.  Susceptible $S(t)$: people before getting infected, but can be infected by neighbors.
2.  Infectious $I(t)$: people who have been infected, and can infect others.

3. Removed $R(t)$: people who endured a complete infection cycle and are immune to the disease.

The number of population is $N(t)$, where $N(t) = S(t) + I(t) + R(t)$. The SIR model describes the change between the three parts of people. We introduce two parameters $\beta$ and $\gamma$.

- $\beta$ Effective contact rate
- $\gamma$ Recovery rate

To be more specific, $\beta$ describes number of susceptible persons that a patient can infect per unit time is proportional to the total number of susceptible persons in the environment. $\gamma$ implies that the number of people removed from the infected person per unit time is proportional to the number of patients. Three differential equations are used to sketch the model:

$$\frac{\mathrm{d}S}{\mathrm{d}t} = -\frac{\beta SI}{N}$$
$$\frac{\mathrm{d}I}{\mathrm{d}t} = \frac{\beta SI}{N} - \gamma I$$
$$\frac{\mathrm{d}R}{\mathrm{d}t} = \gamma I$$

We have to estimate the value of $\beta$ and $\gamma$ simultaneously. The method we use to fit the simple and static SIR is optimization minimizing the sum of squared error of estimated value and real value for infected and recovered population. We use the function provided by python in scipy to do the optimization.

**Graph Neural Netwok**   GNN is a suitable model to disseminate the geographical information as well as the human mobility across location to address the preliminary prediction with a certain amount of factors taken into account. In order to imitate infectious disease modeling, we plan to take create a graph with multiple nodes representing the states of a country and edges defining the spatial and temporal dependencies among states. We present our graph as a multi-layer stack, where each layer manifest the county connectivity graph for one day. Similar as the message-passing framework, we define the update at layer $l$ as

$$m_i^{(l+1)} = \sum_{j \in \mathcal{N}(i)} \mathcal{F}^{(l)}\left(h_i^{(l)}, h_j^{(l)}\right), \quad h_i^{(l+1)} = \mathcal{G}^{(l)}\left(h_i^{(l)}, m_i^{(l+1)}\right)$$

where $\mathcal{F}^{(l)}$ and $\mathcal{G}^{(l)}$ are learned message functions and node update functions respectively, $m^{(l)}$ are the messages passed between nodes, and $h^{(l)}$ are the node representations. Notice that the bottom layer representing the date cases began appearing in the US and the top layer representing the newest date in our dataset.

## Related Works

Various work on SIR model and its extensions has been applied on COVID-19 analysis. Hao et al.[4] reconstruct the outbreak in Wuhan (China) using an extended SAPHIRE model including seven compartments. A CovsirPhy Development Team[5] provided data analysis methods in python package based on SIR-derived models including SIR-F/SIR-FV/SEWIR-F. Yang et al.[6] used previous pandemic data to pretrain the LSTM, and then apply it to predict COVID-19 progression in China. However, different pandemics have different infect ability, so it may lead to inferior prediction results if the model transfer previous pandemic progression directly at the early stage of the pandemic. We will mainly focus on SIR fitting with least squares. We simplify the compartment process since our approach is to apply SIR result in featurization of graph neural network.

Apart from traditional epidemic infectious models like illness incidence forecasting and dynamic disease transmission model as SIR, SEIR[7]. Many novel approaches are based on deep neural networks. Deng et al.[8] proposed a graph message passing framework to combine learned feature embeddings and an attention matrix to model disease propagation over time. Google research team Kapoor et al.[9] apply a simple graph neural network for COVID-19 forecast prediction.

Moreover, some hybrid model has ingeniously combined disease transmission model and gnn model. Jie Zhou, Ganqu Cui etal.[10] combine SEIR and RNN and develop a spatio-temporal graph neural model, where the node feature is based on SEIR and the edge feature is based on the RNN model.

## Datasets and Insights

We make full use of two datasets: the New York Times (NYT) COVID-19 dataset and the Google COVID-19 Aggregated Mobility Research Dataset. The NYT dataset collects all COVID-19 related data and is up to date til Nov.12 2020. It includes but not limits to the number of active cases, confirmed cases and deaths among different locations in the US. The Aggregated Mobility Research Dataset helps us understand the quantity of movement, while the other information combined make distinct node features for NYT. The following figure shows a overview of the infection rate in NY in the United States.

4. Xingjie Hao 2020.
5. Lisphilar 2020.
6. Yang **andothers** 2020.
7. Pei **and** Shaman 2020.
8. Deng **andothers december** 2019.
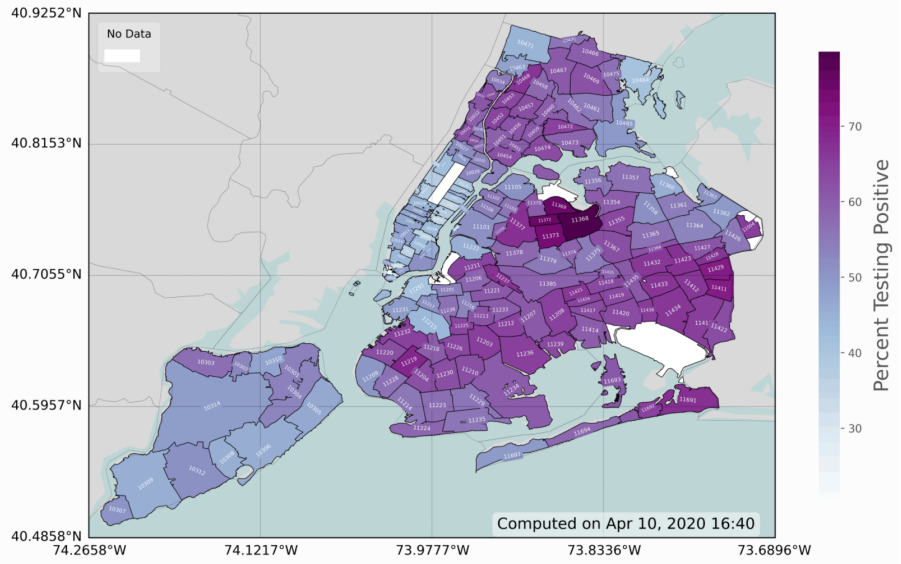9. Kapoor **andothers** 2020.
10. Zheng **andothers** 2020.

**FIGURE 2.** *The zip code map infection rate of New York City*

## Information Visualization

The aim of this part is to provide a initial understanding of the vast amount of data related to COVID-19, based on the insight we gain from the graph, we can embed the essential information to create graph structure from that data and allow users to easily understand node features and edge weights.



(a) Age group difference
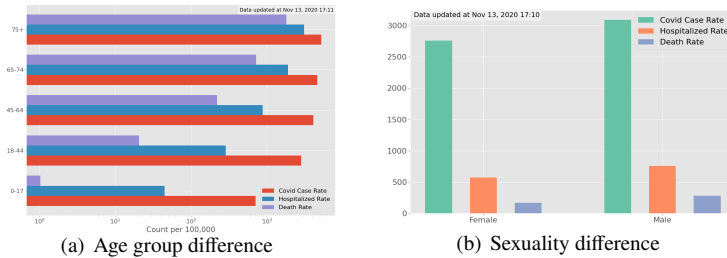
(b) Sexuality difference

**FIGURE 3.** *Covid-19 Cases/Hospitalized/Death rate*

To gain some insights as we are constructing the graph neural network, we notice that different aged groups and different sexuality tend to have different symptoms facing the COVID-19. As a result, we decide to extract the feature and design a factor

representing those information, since the average values across all aged groups or bi-sexuality won't be useful when training our GNN model.

From the above figure, we find that other than the structural node labels, the node information matrix $X$ also provides opportunities to include explicit features, by concatenating each node's embedding vector to its specific row in $X$, we can make the HFSSG model simultaneously learn from all types of features.

## Graph Representation

### Node Embedding

To construct a graph from scratch, we have to define a set of nodes and their edges that denotes connectivity. In case of a geographic map, we take 50 nodes as the states of the United States, and the edges be the common borders between two given states. Therefore, we have 50 nodes in total in our epidemic network. Figure 4 visualizes the embedding of nodes with the name labeling the states:



**FIGURE 4.** *State embedded nodes*

The features of a node is both static and dynamic, we combine both of the two properties. We include 4 variables for the static evaluation, which are population, population density, longitude and latitude. For the dynamic features, we includes the number of cases in the epidemic. The type of cases are confirmed, recovered and death.

### Edge Embedding

For the weight of each edge, we decide to include the properties including the geographical characteristics and population of two connected nodes. We define the

edge weight in the formula,

$$w_{ij} \propto p_i^{\alpha_1} p_j^{\alpha_2} \exp\left(-\frac{d_{ij}}{\alpha_3}\right)$$

where $\alpha_1, \alpha_2, \alpha_3$ are hyper parameters. Through this formula we want to show that the state with larger population are more likely to influence another state connected with it in the network. And the traffic mobility is also an important factor where convenient and developed transports facilitate the propagation of epidemic. And here the distance about traffic we use is the air traffic distance.

## Methods

### SIR with Machine Learning

SIR model divide the population into three compartments, susceptible, infectious and recovered. There are two parameters, $\beta$ and $\gamma$. The problem is how to fitting the two parameters to data. We apply the least square method, which allow us to fit $\beta$ and $\gamma$ simultaneously. It means that the sum of squared errors is a surface in two dimenions and we are looking for the minimum of this surface. This process is referred to as optimization and, fortunately for us, there are many robust algorithms available for this purpose. One of them, Limited Memory BFGS, is implemented in the function `scipy.optimize.minimize` in python. So we can use it to finish the fitting process. Before that, we look at the data of cases in the united states. We focus on the susceptible, infectious and recovered population from the dataset[11] confirmed, recovered and death population.
The plot compares the estimated and real epidemic condition in the USA. It is accurate for the known value, but for estimation, we know that the cases has already exceed the value predicted in the graph. Then we find that the prediction is not that accurate if we want to predict the long-run data based on just a few train set.

### Result Analysis

The reason behind the inaccuracy is that the parameters are always changing with the real events like the release of policies and enhance of people's awareness in wearing masks. So we have to find a way to generate a dynamic or times series method fitting the transmission rate and recovery rate. We have to chop the time serires into multiple time stamps and determine the parameters for each of the period.
Also, the limitation of using only population data is that it omits various information geographically, like the latitude and longitude. It also cannot embed the network structural information between different states like the distance and population density
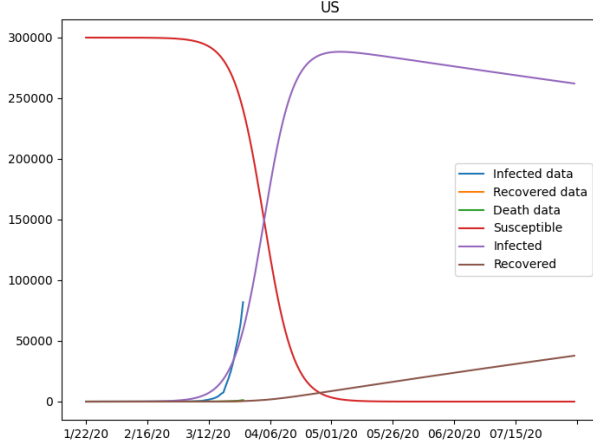
11. Guidotti **and** Ardia 2020.

**FIGURE 5.** *Stable SIR model fitting.*

similarity. Therefore, the original way is not applicable. So we design HFSSG to fit the dynamic parameters.

## Spatio-Temporal Graph Neural Network

Graphs are natural representations for a wide variety of real-life data. Spatio-temporal graphs are a kind of graph that model connections between nodes as a function of time and space, and have found uses in a wide variety of fields.

## Use GCN to Extract Spatial Features

With considering the disease transmission status of similar locations, more accurate prediction can be made. Here we apply the Graph Convolution Networks (GCN) model to extract spatial features.

We construct a two-layer GCN. First we build a sliding window. Let $X_t = [X_1, X_2]$ with $X_1 \in \mathbb{R}^{N \times T \times L_I D_1}$, $X_2 \in \mathbb{R}^{N \times T \times L_I D_2}$, where $N$ is the number of nodes, $T$ is the number of time stamps, $L_1$ is the length of the input sliding window, $D_1$ is the number of static features and $D_2$ is the number of dynamic features. The GCN can be expressed as:

$$z_t = \widehat{A} \, \text{Relu} \left( \widehat{A} X_t W_0 \right) W_1$$

where A denotes the normalized adjacency matrix of the graph G, $W_0, W_1$ denote the weight matrix in the first layer and the second layer.

We use the Graph Attention Networks (GAT) model to extract spatio features. In our model, each node representing one location get information from its neighbours. For example, a city with large amount of population will have more effect on it neighbours. The node feature for similarity can be calculated as:

$$e_{ij} = a\left(W_a \mathbf{z}_t^i, W_a \mathbf{z}_t^j\right)$$

where $a$ is the attention coefficient, $\mathbf{z}_t^i$ is the i-th nodes in $h_t$. $W_a$ is the linear transformation weight matrix which maps the input to a feature space of D-dimension. Then, the final attention coefficient is defined as:

$$a_{ij} = \text{softmax}\left(e_{ij}\right) = \frac{\exp\left(e_{ij}\right)}{\sum_{k=1}^{N} \exp\left(e_{ik}\right)}$$

We use multi-head strategy to calculate **K** independent attention information and sum together to the get representation $\widetilde{\mathbf{z}}_t^i$ for i-th node:

$$\widetilde{\mathbf{z}}_t^i = \sigma\left(\frac{1}{K}\sum_{k=1}^{K}\sum_{j=1}^{N} a_{ij}^k W^k z_t^i\right)$$

Here, $W^k$ is the weight matrix for the k-th head.

**Use RNN to Extract Temporal Features**

The node embedding contains spatial features extracted from the graph. Despite the spatial features, pandemic prediction is also related with temporal features. We also want to use historical temporal patterns to make a better prediction. We input the node embedding into Gate Recurrent Unit (GRU) in order to learn temporal features. We use max-pooling to integrate node embeddings for each location, which is:

$$\widetilde{Z}_t = \text{maxpool}(\widetilde{Z_t^0}, \widetilde{Z_t^1}, ..., \widetilde{Z_t^N})$$

By using the MaxPooling operator, the most important features of all of the nodes are selected, and the dimension of the node information is reduced.

We use Gated Recurrent Unit (GRU) to extract temporal features. GRU is a recurrent neural network that has good performance on temporal sequences modeling. The GRU's hidden representation can be calculated as:

$$h_t = \text{GRU}(\widetilde{Z}_1, \widetilde{Z}_2, ..., \widetilde{Z}_t)$$

For each location, we build a model belonging to that location. The location index $i$ is omitted. $h_t$ therefore contains both spatial and temporal features based on our STGNN model.

**Hybrid Framework using SIR and ST-GNN**

Since we have learned the representation of the nodes in epidemic network. Then we need to combine the STGNN with SIR model, which means using the embedding value to fit the parameters of SIR model.

As we have discussed in the result analysis part of SIR with Machine Learning section, the original model fails to catch the dynamic characters of the epidemic transmission. Therefore, we derive a way to split the entire time into multiple slots. Define prediction window $L_p$ where $\beta$ and $\gamma$ keep constant within one window. A time stamp represents a prediction window, and there are $T/L_p$ sections in total. Each time step is symbolised using index $t$. So our tast is to predict $\beta_t, \gamma_t$ for all the time stamp $t$. We use Multi-Layer Perceptron and sigmoid activation for the neural network construction.

$$\beta_t, \gamma_t = \text{sigmoid}\left(\text{MLP}\left(\boldsymbol{h}_t\right)\right)$$

where MLP indicates Multi-Layer Perceptron, and $h_t$ is the output from the former embedding.

We can also predict the increment of infected and recovered cases, denoted as $\widehat{\Delta I_t}$ and $\widehat{\Delta R_t}$. Similar to the method we use for the parameter fitting, we also use Multi-Layer Perceptron.

$$\widehat{\Delta I_t}, \widehat{\Delta R_t} = \text{MLP}\left(\boldsymbol{h}_t\right)$$

where $h_t$ is the output from the former embedding.

$$\widehat{\boldsymbol{I}_t} = I_{t-1} + \sum\left(\widehat{\Delta I_t}\right)$$
$$\widehat{\boldsymbol{R}_t} = R_{t-1} + \sum\left(\widehat{\Delta \boldsymbol{R}_t}\right)$$

Then we can derive the infected cases and recovery cases through a cumulative sum over the increment value of them.

We define the loss function in the normal way using Mean Squared Error(MSE). We can optimize our estimated value through restricting it close to the real value.

$$L = \left(\widehat{I_t} - I_t\right)^2 + \left(\widehat{\boldsymbol{R}_t} - \boldsymbol{R}_t\right)^2$$

where $\widehat{I_t}$ and $\widehat{\boldsymbol{R}_t}$ denotes the estimated value and $\boldsymbol{I}_t$ and $\boldsymbol{R}_t$.

## Results

Utilizing the $\boldsymbol{I}_t$ and $\boldsymbol{R}_t$ given from the Multi-layer Perceptron, we could calculate and represent the predicted case information as normalized vectors that happened during the last d days. We use `networkx` and `kamada_kawai_layout` scheme to represent the structure of network and death count numbers published by the New York Times,

which includes daily reports of new infections and deaths at both state and county level in US. As is shown in figure 6, the darker the color indicates more cases of COVID-19 locally, which usually also has more edges connected to other nodes.
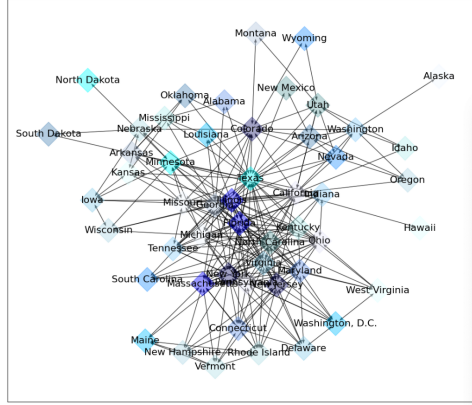


**FIGURE 6.** *Nodes with death counts*

In order to properly utilize the geo-information we embed in the nodes, we locate all nodes on a us map based on its latitude and longitude. For a better visualization, we consider the basic reproduction number as an important factor to denote the contagiousness of a node among its neighbors. Given $d_i$ as the degree of each node, we estimate $q$ as the proportion of its neighbors subsequently infected and average, thus we define the basic reproduction number as

$$R_0 = q * d * \frac{\text{avg}\left(d_i^2\right)}{(\text{avg}\, d_i)^2}$$

Consider for all reproduction values ranging from 0 to 1, we denote the top 10% with higher probability to cause infection cascade as the strong nodes, while for the rest are denoted as weak nodes. For the edges connecting the nodes, we only preserve the top 50% weighted neighbouring edges for simplicity. Figure 7 shows the combined geo-info graph structure map.

**Sensitive Analysis**

HFSSG model is compared with some of the existing baselines:
1. **SIR**: the Susceptible Infectious Removed model. It's a simple and traditional epidemiological model.
2. **SEUR**: the Susceptible Exposed Infectious Removed model. SEIR has one more exposed state compared with the SIR model.
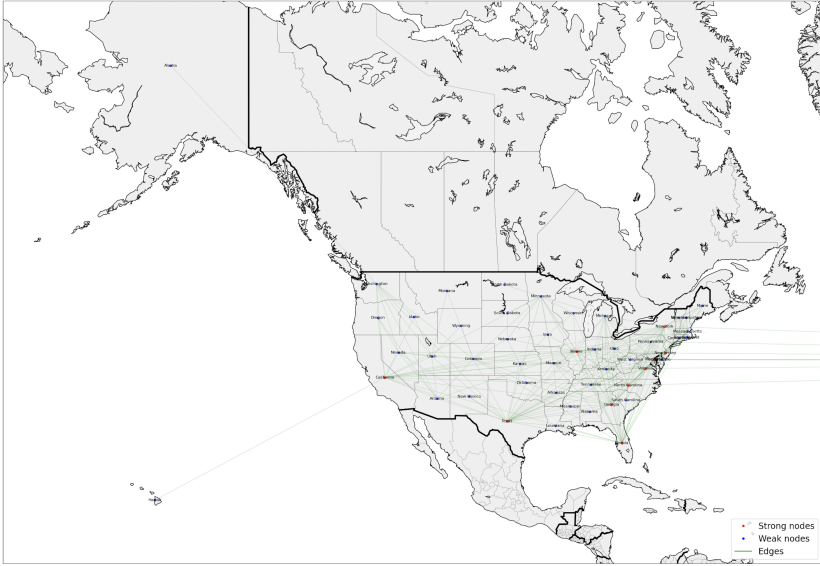3. **HFSSG**: Our proposed model.

**FIGURE 7.** *Graph structured Geo-Info Map*

The prediction is in state-level scale. In order to evaluate the ability of HFSSG, we use the mean square error (MSE) and mean absolute error (MAE). We established a table based on this evaluation method. We use different history data window to compare the performance.

Performance comparison with data window size 5:

| Model | MSE | MAE |
|-------|-----|-----|
| SIR | 2965210 | 910.08 |
| SEIR | 1950414 | 710.34 |
| HFSSG | 421372 | 264.13 |

We also calculated the Performance comparison with data window size 20:

| Model | MSE | MAE |
|-------|-----|-----|
| SIR | 47592364 | 3680.85 |
| SEIR | 28741861 | 2836.37 |
| HFSSG | 5256103 | 1204.36 |

From the results, we can see the HFSSG has a more accurate prediction than SIR and SEIR, both in short time period and long time period.

**Strength and Weaknesses**

We discuss the strength and weaknesses of our model. First from the result section, we can find our model performance exceeds that of SIR and SEIR. The main reason is that we use state-level data and include location information to predict the epidemic propagation.

We consider this dependency between two neighboring states where the communication is more frequently due to the transports convenience. Through network embedding with deep learning methods, we build a clear pipeline from raw information to low-dimensional representations, which allows us to include geographical characters during prediction work. The ability to explain hidden features is why our model performance exceeds that of traditional SIR and SEIR.

In node embedding process, we also consider some other factors like population density and population. They also reflects the similarity between two nodes in the sight of population similarity.

Although the deep learning method achieves greater performance in the prediction work, it is also limited by some weaknesses. One of them exists in the prediction window design. The partition of time series data provides us more flexibility but also arouses some trouble. The data quality is fluctuating during the epidemic process due to multiple reasons like error cases counting due to unfamiliar about the symptom of the disease. Therefore, it could happen that inside some time section the data has low quality, and thus influencing the quality of the prediction work. Moreover, if some emergency policies are announced or significant events take place, resulting in sudden change in the cases number, the prediction inside such window may output less meaningful result.

We can also improve the design for loss function to add more constraints and include more situations in reality. Like what we have referred, different policies lead to change in transmission rate and recovery rate. We generally use MSE to optimize the learning of number of cases and parameters in SIR model. It is possible to add more loss representations for other interpretations.

## Conclusion

In this project report, our proposed HFSSG model is able to make reasonable prediction over traditional model such as SIR and SEIR, which is shown in final values for MSE comparison. Despite the SIR model performs seemingly better with the parameters trained with machine learning, our model is expected to perform better with a longer period of time, since the MLP could aggregate the values from the neighbors and utilize the temporal aspect of our data. Our abstraction of edges assumed that there is a communication between the nodes, which may represent the fact that the flow of people between different states is not being canceled. Clearly due to the quarantine, the edges are shut down and will definitely modify the prediction results like the one presented here. In the future, we plan to tune the hyper-parameters given the prediction

results with the ground truth, and we can identify more representative descriptors for the spread of COVID-19.

# Appendix

## Visualization

### Featurization of sexuality

```python
import csv
import matplotlib.pyplot as plt
import datetime
import numpy as np

plt.style.use('ggplot')

age_data = './boroughs-by-sex.csv'
txt = []

with open(age_data, 'r', newline='') as file:
    reader = csv.reader(file)
    for line in reader:
        txt.append(line)

header = txt[0]

fig,ax = plt.subplots(figsize=(12.5,7))
spacer = -0.25
cii = 0
for plot_indx in range(4,7):
    data_to_plot,x_range = [],[]
    for jj in range(1,len(txt)-1):
        x_range.append(txt[jj][0])
        data_to_plot.append(float(txt[jj][plot_indx]))
    x_plot = np.arange(0,len(x_range))+spacer
    hist =
        ax.bar(x_plot,data_to_plot,label=header[plot_indx].replace('_','
        ').title(),width=0.15,color=plt.cm.Set2(cii))
    spacer+=0.25
    cii+=1

ax.set_xticks(np.arange(0,len(x_range)))
ax.set_xticklabels(x_range)
```

```
33  ax.legend(fontsize=16)
34  ax.tick_params('both',labelsize=16)
35  # textbox showing the date the data was processed
36  txtbox = ax.text(0.0, 0.975, 'Data updated at
    ↪  '+datetime.datetime.now().strftime('%b %d, %Y %H:%M'),
    ↪  transform=ax.transAxes, fontsize=14,
37         verticalalignment='center', bbox=dict(boxstyle='round',
           ↪  facecolor='w',alpha=0.5))
38  txtbox.set_x(0.36-(txtbox.figure.bbox.bounds[2]-(txtbox.clipbox.bounds[2]-
39  fig.savefig(header[0]+'_in_nyc.png',dpi=300,facecolor='#FCFCFC',bbox_inches
    ↪  = 'tight')
40  plt.show()
```

**Featurization of different age group**

```
1   import csv
2   import matplotlib.pyplot as plt
3   import datetime
4   import numpy as np
5
6   plt.style.use('ggplot')
7
8   age_data = './boroughs-by-age.csv'
9   txt = []
10
11  with open(age_data, 'r', newline='') as file:
12      reader = csv.reader(file)
13      for line in reader:
14          txt.append(line)
15
16  header = txt[0]
17
18  fig,ax = plt.subplots(figsize=(20,10))
19  spacer = -0.25
20  for plot_indx in range(4,7): ## Only focus on
    ↪  BK_CASE_RATE,BK_HOSPITALIZED_RATE,BK_DEATH_RATE
21      data_to_plot,x_range = [],[]
22      for jj in range(1,len(txt)-1):
23          x_range.append(txt[jj][0])
24          data_to_plot.append(float(txt[jj][plot_indx]))
25      # print(data_to_plot)
26      x_plot = np.arange(0,len(x_range))+spacer
```

```
27      hist =
        ↪ ax.barh(x_plot,data_to_plot,label=header[plot_indx].replace('_',
        ↪ ').title(),height=0.25,log=True)
28      spacer+=0.25
29
30 ax.set_xlabel('Count per 100,000',fontsize=20)
31 ax.set_yticks(np.arange(0,len(x_range)))
32 ax.set_yticklabels(x_range)
33 ax.legend(fontsize=16)
34 ax.tick_params('both',labelsize=16)
35 # fig.suptitle('COVID-19 in NYC by '+header[0].replace('_','
   ↪ ').title(),x=0.4,y=0.92,fontsize=18)
36
37 txtbox = ax.text(0.0, 0.975, 'Data updated at
   ↪ '+datetime.datetime.now().strftime('%b %d, %Y %H:%M'),
   ↪ transform=ax.transAxes, fontsize=14,
38        verticalalignment='center', bbox=dict(boxstyle='round',
          ↪ facecolor='w',alpha=0.5))
39 txtbox.set_x(1.04-(txtbox.figure.bbox.bounds[2]-(txtbox.clipbox.bounds[2]
40 fig.savefig(header[0]+'_in_nyc.png',dpi=300,facecolor='#FCFCFC',bbox_incl
   ↪ = 'tight')
41 plt.show()
```

**Graph Structured Representation**

```
1 from torch_geometric.utils.convert import to_networkx
2 import matplotlib.pyplot as plt
3 import networkx as nx
4 from src.graph_loader import make_fc_graph
5 from pathlib import Path
6 from src.data_loader import load_pickle
7 import pickle
8 from mpl_toolkits.basemap import Basemap as Basemap
9 import pandas as pd
10 import matplotlib.lines as mlines
11
12 SOURCE_PATH = Path(__file__).parent
13
14 def visualize_on_map(data, labels):
15      df = pd.read_csv(SOURCE_PATH /
          ↪ '../dataset/lat_lon_usa.csv')
16
17      plt.figure(figsize=(300, 400))
18      m = Basemap(
```

```python
                  projection='merc',
                  llcrnrlon=-180,
                  llcrnrlat=10,
                  urcrnrlon=-50,
                  urcrnrlat=70,
                  lat_ts=0,
                  resolution='l',
                  suppress_ticks=True)

        data_graph = to_networkx(data)
        pos = {}
        temp_x, temp_y = [], []
        edge_weights = data.edge_attr.numpy()
        degrees = dict(data_graph.degree)
        edges = data_graph.edges
        strong_nodes = []
        weak_nodes = []
        node_strengths = {}
        weightage = 0.5
        for index, edge in enumerate(edges):
                n1, n2 = edge

                if n1 in node_strengths:
                        node_strengths[n1] +=
                        ↪ edge_weights[index]
                else:
                        node_strengths[n1] = edge_weights[index]
                        ↪ + weightage*degrees[n1]
                if n2 in node_strengths:
                        node_strengths[n2] +=
                        ↪ edge_weights[index]
                else:
                        node_strengths[n2] = edge_weights[index]
                        ↪ + weightage*degrees[n2]

        print(degrees)
        print(labels)
        for k, v in node_strengths.items():
                print("State: ", labels[k], " Strength: ", v)
        # The top x nodes get assigned the strong nodes
        ↪ category
        x = 10
        counter = 0
```

```
57          for k, v in sorted(node_strengths.items(), key=lambda
        ↪   item: item[1], reverse=True):
58              if counter < x:
59                      strong_nodes.append(k)
60              else:
61                      weak_nodes.append(k)
62
63              counter += 1
64
65          # Assign cor-ordinates to plot the graph
66          for index, node in enumerate(data_graph.nodes):
67              state = labels[node]
68
69              lon = df.loc[df['State'] == state, 'lon'].item()
70              lat = df.loc[df['State'] == state, 'lat'].item()
71              temp_x.append(lon)
72              temp_y.append(lat)
73          mx, my = m(temp_x, temp_y)
74
75          for index, val in enumerate(data_graph.nodes):
76              pos[val] = (mx[index], my[index])
77
78          nx.draw_networkx_nodes(data_graph, pos=pos,
        ↪   nodelist=strong_nodes, node_size=20, node_color='r',
        ↪   alpha=0.8)
79          nx.draw_networkx_nodes(data_graph, pos=pos,
        ↪   nodelist=weak_nodes, node_size=15, node_color='b',
        ↪   alpha=0.6)
80          nx.draw_networkx_edges(data_graph, pos=pos,
        ↪   edgelist=data_graph.edges, width=1.0, alpha=0.1,
        ↪   arrows=False, edge_color='g')
81
82          nx.draw_networkx_labels(data_graph, pos=pos,
        ↪   labels=labels, font_size=8)
83
84          m.drawcountries(linewidth=3)
85          m.drawstates(linewidth=0.2)
86          m.drawcoastlines(linewidth=1)
87          m.fillcontinents(alpha=0.3)
88          line1 = mlines.Line2D(range(1), range(1), color="white",
        ↪   marker='o', markerfacecolor="red")
89          line2 = mlines.Line2D(range(1), range(1), color="white",
        ↪   marker='o', markerfacecolor="blue")
```

```python
90          line3 = mlines.Line2D(range(1), range(1), color="green",
        ↪ marker='', markerfacecolor="green")
91          plt.legend((line1, line2, line3), ('Strong nodes', 'Weak
        ↪ nodes', 'Edges'),
92                          loc=4, fontsize='xx-large')
93          plt.savefig(SOURCE_PATH / '../results/usmap.png')
94          plt.show()
95
96
97  def visualize_graph(data, labels, positions=None):
98          colors = [i for i in range(17)]*3
99          colors = [0.00001*i for i in colors]
100         data_graph = to_networkx(data)
101         pos = nx.kamada_kawai_layout(data_graph)
102
103         plt.figure(1, figsize=(14, 12))
104
105         nx.draw_networkx_nodes(data_graph, pos=pos,
        ↪ nodelist=data_graph.nodes,
        ↪ alpha=0.1,node_size=700,cmap=plt.cm.jet)
106         # nx.draw(data_graph, cmap=plt.get_cmap('Set1'),
        ↪ node_size=75, linewidths=6)
107         nx.draw_networkx_edges(data_graph, pos=pos,
        ↪ edgelist=data_graph.edges, width=1.0, alpha=0.5)
108
109         nx.draw_networkx_labels(data_graph, pos=pos,
        ↪ labels=labels, font_size=16)
110         # plt.savefig(SOURCE_PATH / '../results/usnode2.png')
111         plt.show()
112
113
114  if __name__ == '__main__':
115         metadata, data = make_fc_graph()
116
117         with open(SOURCE_PATH /
        ↪ '../dataset/generated/usa/node_to_ss', 'rb') as
        ↪ file:
118                 node_to_ss = pickle.load(file)
119
120         visualize_graph(data, node_to_ss)
121         # visualize_on_map(data, node_to_ss)
122         # Load to node to ss dict
```

**Traditional SIR Method**

```
1  library (readr)
2  library("COVID19")
3  require(deSolve)
4  usa_data <- covid19('USA')
5  par(mfrow=c(1,4))
6  plot(x = usa_data$date, y = 326687501 - usa_data$deaths -
   ↪ usa_data$confirmed - usa_data$recovered, type = 'l', col =
   ↪ 4,xlab = "Time/month", ylab = 'Population', main =
   ↪ 'Susceptible')
7  plot(x = usa_data$date, y = usa_data$confirmed -
   ↪ usa_data$recovered - usa_data$deaths, type = 'l',col = 2,
   ↪ xlab = "Time/month", ylab = 'Population', main = 'Infected')
8  plot(x = usa_data$date, y = usa_data$recovered, type = 'l',col =
   ↪ 3, xlab = "Time/month", ylab = 'Population', main =
   ↪ 'Recovered')
9
10 plot(x = usa_data$date, y = 326687501 - usa_data$deaths -
   ↪ usa_data$confirmed - usa_data$recovered,
11     type = 'l',col = 2, xlab = "Time/month", ylab =
       ↪ 'Population', main = 'COVID-19 case',
12     ylim = c(0, 3.3*10E7)
13     ,yaxs = "i"
14     )
15 lines(x = usa_data$date, y = usa_data$confirmed -
   ↪ usa_data$recovered - usa_data$deaths, col = 4)
16 lines(x = usa_data$date, y = usa_data$recovered, col = 3)
17 legend("left", legend=c("Susceptible", "Infected","Recovered"),
18        col=c("red", "blue", "green"), lty=1)
```

**Training ST-GNN**

```
1  import numpy as np
2  import torch
3  import torch.nn as nn
4  import torch.nn.functional as F
5  from utils import Data2Graph
6  import scipy.stats
7  from dgl import DGLGraph
8  import sklearn
9  import model
10
```

```
11   device = torch.device("cuda:0" if torch.cuda.is_available()
     ↪ else "cpu")
12
13   data_filename='FILEPATH'
14   auxiliary_filename='uszips.csv'
15
16   slot_len = 6
17   pred_horizon = 60
18   step_size = 5
19   mid_time = (pred_horizon//2)+1
20
21   nhid1 = 400
22   nhid2 = 200
23   gru_dim = 100
24   num_heads = 3
25   lr = 1e-3
26   r = 1
27
28   num_epochs = 2000
29
30
31   feat_tensor, Adj, active_cases, confirmed_cases, popn,
     ↪ selected_counties,g = Data2Graph.load_data(data_filename,
     ↪ auxiliary_filename,active_thresh=1000,r=r)
32   recovered_cases = confirmed_cases - active_cases
33   state_names=np.array(selected_counties['state'])
34
35   d_active = []
36   d_recovered = []
37   for i in range(active_cases.size(0)):
38     i_active = [0]
39     i_recovered = [0]
40     for j in range(1,active_cases.size(1)):
41       i_active.append(active_cases[i,j]-active_cases[i,j-1])
42
         ↪ i_recovered.append(recovered_cases[i,j]-recovered_cases[i,j-1])
43     d_active.append(i_active)
44     d_recovered.append(i_recovered)
45
46   d_active = torch.tensor(d_active).to(device)
47   d_recovered = torch.tensor(d_recovered).to(device)
48
49
```

```
50   N,T,D = feat_tensor.shape
51   nfeat = slot_len*D
52   nclass = pred_horizon*3
53   attn_dim=int(nfeat*1.5)
54
55
56   num_pred_time = len(range(slot_len,T))
57
58   num_val_pred_time = 60
59   num_test_pred_time = 5
60   num_train_pred_time = num_pred_time-num_val_pred_time -
     ↪   num_test_pred_time
61
62   test_pred_time=np.arange(T-num_test_pred_time,T)
63   val_pred_time=np.arange(test_pred_time[0]-num_val_pred_time,test_pred_tin
64   train_pred_time =
     ↪   np.arange(slot_len,val_pred_time[0]+pred_horizon-1)
65
66   train_input_time = np.arange(train_pred_time[-1]+1-pred_horizon)
67   val_input_time =
     ↪   np.arange(val_pred_time[0]-slot_len,val_pred_time[-1]+1-pred_horizon)
68   test_input_time = np.array([78,79,80,81,82,83])
69
70
71   # train, val, test split along the spatial dimension
72
73   num_train = int(1*N)
74   num_val = int(0*N)
75   num_test=N-num_train-num_val
76
77   #shuffle_order = np.random.permutation(N)
78   shuffle_order=np.arange(N)
79
80   train_locs = np.arange(N)[shuffle_order[:num_train]]
81   val_locs =
     ↪   np.arange(N)[shuffle_order[num_train:num_train+num_val]]
82   test_locs = np.arange(N)[shuffle_order[num_train+num_val:]]
83
84
85   # data splitting
86
87   W_train=Adj[train_locs,:][:,train_locs].to(device)
88   train_features=feat_tensor[train_locs,:,:][:,train_input_time,:][:].to(de
```

```
89   train_labels_active =
     ↪   active_cases[train_locs,:][:,train_pred_time].to(device)
90   train_d_active =
     ↪   d_active[train_locs,:][:,train_pred_time].to(device)
91   train_labels_confirmed =
     ↪   confirmed_cases[train_locs,:][:,train_pred_time].to(device)
92   train_labels_recovered =
     ↪   recovered_cases[train_locs,:][:,train_pred_time].to(device)
93   train_d_recovered =
     ↪   d_recovered[train_locs,:][:,train_pred_time].to(device)
94   train_popn=popn[train_locs,:][:].to(device)
95

96

97   # For validation and testing, we use the full training data
     ↪   points as well

98

99   val_set = np.concatenate((train_locs, val_locs),axis=0)

100

101  W_val=Adj[val_set,:][:,val_set].to(device)
102  val_features = feat_tensor[val_set,:,:][:,
     ↪   val_input_time,:][:].to(device)
103  val_labels_active =
     ↪   active_cases[val_set,:][:,val_pred_time].to(device)
104  val_d_active = d_active[val_set,:][:,val_pred_time].to(device)
105  val_labels_confirmed =
     ↪   confirmed_cases[val_set,:][:,val_pred_time].to(device)
106  val_popn=popn[val_set,:][:].to(device)
107  val_d_recovered =
     ↪   d_recovered[val_set,:][:,val_pred_time].to(device)

108

109  # For testing, we use the full training set as well for graph
     ↪   based inference

110

111  test_set = np.concatenate((train_locs, test_locs),axis=0)

112

113  W_test=Adj[test_set,:][:,test_set].to(device)
114  test_features = feat_tensor[test_set,
     ↪   :,:][:,test_input_time].to(device)
115  test_labels_active =
     ↪   d_active[test_set,:][:,test_pred_time].to(device)
116  test_d_active =
     ↪   active_cases[test_set,:][:,test_pred_time].to(device)
```

```
117   test_labels_confirmed =
      ↪ confirmed_cases[test_set,:][:,test_pred_time].to(device)
118   test_d_recovered =
      ↪ d_recovered[test_set,:][:,test_pred_time].to(device)
119
120
121   for target_loc in range(45):
122     print('----------Start training for loc
        ↪ %d----------'%target_loc)
123     num_train_steps = len(train_input_time)-slot_len+1
124     batch_feat = []
125     batch_active = []
126     batch_recover = []
127     batch_I = []
128     batch_R = []
129     batch_S = []
130     batch_It = []
131     batch_Rt = []
132
133     for t in range(0, num_train_steps, step_size):
134       t_feat =
          ↪ train_features[:,t:t+slot_len,:].view(int(num_train),slot_len*D)
135
          ↪ t_active=train_d_active[:][target_loc,t:t+pred_horizon].float()
136       t_recovered =
          ↪ train_d_recovered[:][target_loc,t:t+pred_horizon].float()
137       last_I =
          ↪ active_cases[target_loc,train_pred_time[t]-1].unsqueeze(-1).float
138       last_R =
          ↪ recovered_cases[target_loc,train_pred_time[t]-1].unsqueeze(-1).fl
139       last_S = popn[target_loc,0].to(device).float().unsqueeze(-1)
          ↪ - last_I - last_R
140
141       batch_It.append(d_active[target_loc,
          ↪ train_pred_time[t]-1].float())
142       batch_Rt.append(d_recovered[target_loc,
          ↪ train_pred_time[t]-1].float())
143
144       batch_feat.append(t_feat)
145       batch_active.append(t_active)
146       batch_recover.append(t_recovered)
147       batch_I.append(last_I)
148       batch_R.append(last_R)
```

```
149        batch_S.append(last_S)
150
151
152     batch_It = torch.stack(batch_It).to(device).squeeze()
153     batch_Rt = torch.stack(batch_Rt).to(device).squeeze()
154     batch_feat = torch.stack(batch_feat).to(device)
155     batch_active = torch.stack(batch_active).to(device)
156     batch_recover = torch.stack(batch_recover).to(device)
157     batch_I = torch.stack(batch_I).to(device).squeeze()
158     batch_R = torch.stack(batch_R).to(device).squeeze()
159     batch_S = torch.stack(batch_S).to(device).squeeze()
160
161     valid_feat = []
162     valid_active = []
163     valid_recover = []
164     valid_I = []
165     valid_R = []
166     valid_S = []
167     valid_It = []
168     valid_Rt = []
169
170     num_val_steps = val_features.shape[1]-slot_len+1
171     for t in range(0, num_val_steps, step_size):
172       t_feat =
          ↪  val_features[:,t:t+slot_len,:].view(int(num_train),slot_len*D).floa
173
          ↪  t_active=val_d_active[:][target_loc,t:t+pred_horizon].float()
174       t_recovered =
          ↪  val_d_recovered[:][target_loc,t:t+pred_horizon].float()
175
176       last_I =
          ↪  active_cases[target_loc,val_pred_time[t]-1].unsqueeze(-1).float().t
177       last_R =
          ↪  recovered_cases[target_loc,val_pred_time[t]-1].unsqueeze(-1).float(
178       last_S = popn[target_loc,0].to(device).float().unsqueeze(-1)
          ↪  - last_I - last_R
179
180       valid_It.append(d_active[target_loc,
          ↪  val_pred_time[t]-1].float())
181       valid_Rt.append(d_recovered[target_loc,
          ↪  val_pred_time[t]-1].float())
182
183       valid_feat.append(t_feat)
```

```
184       valid_active.append(t_active)
185       valid_recover.append(t_recovered)
186       valid_I.append(last_I)
187       valid_R.append(last_R)
188       valid_S.append(last_S)
189
190     valid_It = torch.stack(valid_It).to(device).flatten()
191     valid_Rt = torch.stack(valid_Rt).to(device).flatten()
192     valid_feat = torch.stack(valid_feat).to(device)
193     valid_active = torch.stack(valid_active).to(device)
194     valid_I = torch.stack(valid_I).to(device).flatten()
195     valid_R = torch.stack(valid_R).to(device).flatten()
196     valid_S = torch.stack(valid_S).to(device).flatten()
197
198
199
200     # training the model
201     model_name = state_names[target_loc]
202     model=model.GAT(g, nfeat, nhid1, nhid2, gru_dim, N, num_heads,
        ↪  pred_horizon,attn_dim).to(device)
203     model = model.float()
204     optimizer = torch.optim.Adam(model.parameters(), lr=lr)
205
206     criterion = nn.MSELoss(reduction='sum')
207
208     for epoch in range(num_epochs):
209
210       model.train()
211       optimizer.zero_grad()
212
213       active_pred, recovered_pred, phy_active, phy_recover =
        ↪  model(batch_feat, popn[target_loc,0].to(device).float(),
        ↪  batch_I, batch_R,batch_S,batch_It,batch_Rt) # forward
        ↪  pass
214
215       loss =
        ↪  F.mse_loss(active_pred,batch_active)+F.mse_loss(recovered_pred,ba
        ↪  batch_active.flatten())+F.mse_loss(phy_recover,
        ↪  batch_recover.flatten())
216
217       loss.backward()
218       optimizer.step()
219
```

```
220    #Evaluate
221    model.eval()
222    with torch.no_grad():
223      active_val, recovered_val, vphy_active, vphy_recover =
       ↪ model(valid_feat,
       ↪ popn[target_loc,0].to(device).float(), valid_I,
       ↪ valid_R,valid_S,valid_It,valid_Rt) # forward pass
224      mse =
       ↪ sklearn.metrics.mean_squared_error(valid_active.cpu().detach().f
225      ccc =
       ↪ model.ccc(valid_active.squeeze().cpu().detach().flatten().numpy()
226      r2 =
       ↪ sklearn.metrics.r2_score(valid_active.cpu().detach().flatten().nu
       ↪ active_val.cpu().detach().flatten().numpy())
227
228    if epoch==0:
229      min_loss = mse
230      min_r2 = r2
231      torch.save(model.state_dict(),model_name)
232    if mse < min_loss:
233      min_loss = mse
234      min_r2 = r2
235      torch.save(model.state_dict(),model_name)
236    if epoch % 100 == 0:
237      print("Epoch %d, Train loss %f, Val loss %f, Val R2 %f,
       ↪ Val ccc %f, %.2f, %.2f"%(epoch, loss, mse,
       ↪ np.mean(r2), np.mean(ccc), model.alpha_list,
       ↪ model.beta_list))
```

# References

Deng, Songgaojun, Shusen Wang, Huzefa Rangwala, Lijing Wang **and** Yue Ning. **december** 2019. Graph Message Passing with Cross-location Attentions for Long-term ILI Prediction ().

Guidotti **and** Ardia. 2020. Covid-19 Data Hub. Available at <https://covid19datahub.io/>.

Kapoor, Amol, Xue Ben, Luyang Liu, Bryan Perozzi, Matt Barnes, Martin Blais **and** Shawn O'Banion. 2020. Examining COVID-19 Forecasting using Spatio-Temporal Graph Neural Networks. arXiv: 2007.03113 [`cs.LG`].

Lisphilar. 2020. CPython package for COVID-19 analysis with SIR-derived ODE models. Available at <https://github.com/lisphilar/covid19-sir>.

Pei, Sen, **and** Jeffrey Shaman. 2020. Initial Simulation of SARS-CoV2 Spread and Intervention Effects in the Continental US. *medRxiv,* https://doi.org/10.1101/2020.03.21.20040303. eprint: https://www.medrxiv.org/content/early/2020/03/27/2020.03.21.20040303.full.pdf. Available at <https://www.medrxiv.org/content/early/2020/03/27/2020.03.21.20040303>.

Xingjie Hao, et al. 2020. Reconstruction of the full transmission dynamics of COVID-19 in Wuhan. *Nature* 584:420–424.

Yang, Zifeng, Zhiqi Zeng, Ke Wang, Sook-San Wong, Wenhua Liang, Mark Zanin, Peng Liu **andothers**. 2020. Modified SEIR and AI prediction of the epidemics trend of COVID-19 in China under public health interventions. *Journal of Thoracic Disease* 12 (3). ISSN: 2077-6624. Available at <http://jtd.amegroups.com/article/view/36385>.

Zheng, Yunling, Zhijian Li, Jack Xin **and** Guofa Zhou. 2020. A Spatial-Temporal Graph Based Hybrid Infectious Disease Model with Application to COVID-19. *arXiv preprint arXiv:2010.09077*.