

A Microscopic Pandemic Simulator for Pandemic Prediction Using Scalable Million-Agent Reinforcement Learning

Zhenggang Tang¹, Kai Yan^{1*}, Liting Sun², Wei Zhan² and Changliu Liu³

¹Peking University

²University of California, Berkeley

³Carnegie Mellon University

{tangzhenggang, kaiyan}@pku.edu.cn, {litingsun, wzhan}@berkeley.edu, cliu6@andrew.cmu.edu

Abstract

Microscopic epidemic models are powerful tools for government policy makers to predict and simulate epidemic outbreaks, which can capture the impact of individual behaviors on the macroscopic phenomenon. However, existing models only consider simple rule-based individual behaviors, limiting their applicability. This paper proposes a deep-reinforcement-learning-powered microscopic model named Microscopic Pandemic Simulator (MPS). By replacing rule-based agents with rational agents whose behaviors are driven to maximize rewards, the MPS provides a better approximation of real world dynamics. To efficiently simulate with massive amounts of agents in MPS, we propose Scalable Million-Agent DQN (SMADQN). The MPS allows us to efficiently evaluate the impact of different government strategies. This paper first calibrates the MPS against real-world data in Allegheny, US, then demonstratively evaluates two government strategies: information disclosure and quarantine. The results validate the effectiveness of the proposed method. As a broad impact, this paper provides novel insights for the application of DRL in large scale agent-based networks such as economic and social networks.

1 Introduction

A good epidemic prediction model is indispensable to mitigate pandemics such as COVID-19, which could help governments derive optimal policies that balance public health and economic resiliency. Various epidemic models have been proposed in literature, among which microscopic models [Aleta *et al.*, 2020] are particularly useful since they are fine-grained enough to encode individual behaviors and local contact traces. However, in the existing works, individuals are either modeled with fixed behaviors [Eilersen and Sneppen, 2020] or are adhere to a given script or rule [Meloni *et al.*, 2011]. These rules may be too simple to fully capture people’s diverse behaviors under environmental changes. On the other hand, reinforcement learning (RL) has been proved

empirically good at generating complex behaviors, where the agent aims to optimize a relatively simple reward function without being guided by handcrafted scripts. RL models are more explainable and natural, because real-world individuals are driven by different motivations, which is coherent with the reward optimization process. Many recent works apply deep RL on microscopic epidemic models [Kompella *et al.*, 2020]. However, these works mainly focus on policy optimization for governments to balance public health and economy without detailed modeling of individual behaviors.

To address the individual behavior modeling problem in pandemic modelling, this paper proposes Microscopic Pandemic Simulator (MPS), a novel microscopic epidemic model where individuals are controlled by multi-agent (MA) RL policies and thus able to change their behaviours according to information gained from their surroundings and the government, as inspired by and expanded upon Liu’s previous work[Liu, 2020]. Since most practical epidemic models contain millions of agents [Hoertel *et al.*, 2020], the main challenge of applying MARL is its scalability. Without special design, it is impossible to make standard RL algorithms applicable to such a massive model. While there are previous works on MARL with million-level number of agents [Zheng *et al.*, 2017], the pandemic environment imposes a much harder challenge for two reasons: 1) rewards of actions are significantly delayed up to several days due to the presence of incubation period, causing 1-step TD learning improper; 2) Our problem is non-ergodic: agents’ joint-state significantly changes as the epidemic spreads, aggravating oscillation during training. To solve these problems with such a large scale, we proposed Scalable Million-Agent DQN (SMADQN), a novel DQN-based algorithm with specially designed replay-buffer and processes for calculating TD(λ), which solves the difficulties above well with such scalability.

To validate authenticity and adaptability of our model, we apply our model on a large scale COVID-19 simulation. Unlike existing works that consider simple contact networks with limited numbers of facilities or complex networks within a city [Aleta *et al.*, 2020], this paper provides a large-scale county-level simulation with all population in the Allegheny county. It includes more than 10^6 residents and a comprehensive contact network including 14 types of facilities. By building a novel demographic dataset of Allegheny County with contact network, we show that the model fit real-world

*Equal contribution for the first two authors.

data well. We further tested the performance of two progressive government strategies on our model, namely, *Information Disclosure* (ID) and *Quarantine* (QT). ID means that the government will disclose information of the pandemic to residents, such as the number of infections in each facility. QT refers to the extra government mitigation strategy that requires symptomatic individuals and part of close contacts of them to be quarantined. We conclude that ID can help mitigate the epidemic and reduce people’s activity levels. However, it is not strong enough to completely control the epidemic. Meanwhile, QT, although posing higher requirements for governments’ administrative capability, can control the epidemic with minimum negative impacts on economic and social activities.

Our contributions are two-fold. The **contributions to epidemic modeling** are the following. 1) We proposed the Microscopic Pandemic Simulator (MPS) where individuals as modeled as rational agents instead of following simple rules. 2) We built a comprehensive demographic dataset with a contact network of Allegheny County. To the authors’ best knowledge, this is the first work in such details in this area. 3) With the MPS and the dataset, we have provided a flexible tool to test the impacts of different governmental strategies. The **contribution to MARL** is the novel soft-DQN algorithm, i.e., SMADQN. It improves the efficiency of learning by utilizing λ -return for million-level multi-agent problem with delayed reward signal and ever-changing joint-state.

2 Related Work

Microscopic Epidemic models. Microscopic epidemic models contain three parts [Hoertel *et al.*, 2020]: personal status, contact network between individuals and a reasonable synthetic population. Generally, personal status includes health status, vulnerability data (e.g., age, basic medical condition [Chang *et al.*, 2020]) and social contacts; contact network is usually divided into different layers, each of which consists of a graph with clusters, representing different places (e.g., workplace and household); synthetic population is retrieved by either mobility data such as GPS [Aleta *et al.*, 2020] or generated from a certain distribution [Eilersen and Sneppen, 2020; Silva *et al.*, 2020] coherent to census data. To balance the granularity and accuracy against computational resources and data available, our work adopts basic personal status settings, fine contact network modeling with 14 layers, and generated synthetic population merged from ArcGIS [ArcGIS, 2020] and US synthetic population dataset [Wheaton, 2012].

Large-Scale Multi-Agent RL. Our problem is framed as a large-scale MA problem. Qu and Li [Qu and Li, 2019] [Qu *et al.*, 2020] considered cooperative games and exploited problem-specific structures, such as a tree or a network. In our work, agents also interact in a network. Instead of focusing on policy optimization, we exploit the explicitness of reward parameters and inherent ability of adaption of RL to produce a more explainable and flexible model in a non-cooperative setting. The work that proposed the most closely related algorithm to ours is [Yang *et al.*, 2018], where they simulated a million-level prey-predator world and proposed a DQN algorithm with redesigned replay buffers. However,

their environment has an infinite horizon, and is ergodic: all agent’s joint-states in different steps are similar, which makes the RL algorithm easier to converge. Furthermore, their reward signals have no delay, while our SMADQN algorithm is designed to solve the non-ergodic environment with long episodic length and delayed reward signal.

3 Problem Formulation

In this work, we model epidemics using Multi-Agent Partial Observation Markov Decision Process (MA-POMDP) where every susceptible individual is viewed as one agent. We are interested in systems that contain more than $n > 10^6$ agents. The MA-POMDP is defined by $(\mathcal{S}, \mathcal{O}, \mathcal{A}, R, P)$. \mathcal{S} represents the joint-state space, and $\mathcal{O} = \otimes_{i \in \{1, 2, \dots, n\}} \mathcal{O}_i$ is the union of all n agents’ observation spaces. \mathcal{A} is the action space for each agent. Agent i ’s action is denoted a_i . $R(s, a_1, a_2, \dots, a_n; i)$ is the reward function for agent i . $P(s' | s, a_1, a_2, \dots, a_n)$ is transition probability from state s to state s' when agent i takes action a_i for all i . Agent i has a stochastic policy $\pi_i(o_i, a_i; \theta_i)$ parameterized by θ_i which outputs the probability to take action a_i when observing o_i . Every agent i optimizes its expected accumulative reward $U_i(\theta_i) = E_{a_1 \sim \pi_1, a_2 \sim \pi_2, \dots, a_n \sim \pi_n, [\sum_t \gamma^t R(s^t, a_1^t, a_2^t, \dots, a_n^t; i)]}$ with the discounted factor γ , where t indicates discrete time steps.

Actions: Agent i ’s action is modeled into three parts: the activity level, whether to wear a mask, and shopping decision, i.e., $a_i = a_{i,act} \times a_{i,mask} \times a_{i,shop}$. To avoid being infected, the most effective way is staying away from risky (with high probability of infection) and non-compulsory facilities. However, it is cumbersome and requires a large action space to allow the agent to choose which facility to visit. So, we assign each facility with a risk level: $MinAct$, and only enable the agent to decide a discrete activity level, then let it visit facilities whose risk level, $MinAct$, are under this level. Note that $MinAct$ of workplaces, schools and households are all set to 0 as they are compulsory although visiting these facilities may also be risky. That’s how we model $a_{i,act}$ and it’s still consistent with reality since it’s natural for individual to stop visiting the riskiest facility first for health. Besides activity level, wearing reduces infection probability during contact while also brings with a little inconvenience, and we build its decision as $a_{i,mask}$. Furthermore, we separate the decision of shopping $a_{i,shop}$, which indicates whether to go shopping in retail stores, or shopping online, or not shopping at all, from a_{act} because $a_{i,shop}$ is influenced by both the epidemic and supply level of home.

Observations: Agent i ’s default observation is $o_{i,no\ info} = o_{i,hea} \times o_{i,hou} \times o_{i,sup} \times o_{city} \times o_t$. Like any real-life individual, an agent gathers information from both its surroundings and the government. For the former part, we assume that they always know the current health condition of itself ($o_{i,hea}$) and other people living in the same household ($o_{i,hou}$), as well as the amount of necessary stocks for living ($o_{i,sup}$), such as food; for the latter, we choose the total number of infection (o_{city}) and the number of days since the first cases are discovered locally (o_t), which are accessible from almost every government in real life. With information disclosure, the observation space will be

increased by one dimension, i.e., $o_{i,sur}$, which indicates the severity of infections in each facility that agent i may visit calculated by the government.

States: Similar to the observation, $S = \{S_i, i \in \{1, 2, \dots, n\}\} \times s_t$ and $S_i = s_{i,hea} \times s_{i,sup}$. The only difference between state and observation is health condition, as people cannot tell whether they are infected under the presence of incubation period.

Rewards: Reward of an agent reflects the incentive of balancing the two goals: ‘‘avoiding infection’’ and ‘‘maintaining normal life’’. The former incentive results in a huge one-time negative reward R_{ill} when getting infected (smoothed in practice; see appendix for details) and a $R_{shop} = -1$ for offline shopping. The latter leads to 3 components of reward: a positive reward $R_{act}(a_i)$ increased with higher activity level $a_{i,act}$; a negative reward R_{mask} for wearing masks, as many people are reluctant to wear masks; and a negative reward R_{sup} as the stocked supplies (e.g., food) decreased with lower supply level $o_{i,sup}$. To further adapt to real-world situation, we impose a more negative reward R_{eth} , the ethical penalty, for higher activity levels or not wearing mask when having symptoms. This is because real-world individuals are not fully selfish: they will avoid infecting others and comply with the restrictions once fallen ill.

More details of the observation space, action space, reward and other RL parameters can be found in the appendix. With these definitions, we will derive the transition function of the simulator in the following sections.

4 Microscopic Pandemic Simulator

This section introduces the basic settings of Microscopic Pandemic Simulator (MPS). MPS takes one day as a discrete round, which is the finest practical grain of time since hour-level simulation requires prohibitively expensive computational resources and GPS data for commute patterns are absent. We show in our experiment that this grain of time is fine enough for authentic simulation. Figure 1 provides an overview of the pipeline of the simulator and its interaction with the SMADQN-controlled agents. The pipeline is organized in a logical order: At the beginning of each day, the government implements some mitigation strategy (e.g., adjust maximum capacity ratio for restaurants); then, agents will observe the current situation and decide whether to visit different facilities, to wear masks and to shop for this day. With determined visit list for each facility, the disease will spread from the infected agents to healthy ones within each facility. Finally, modified by the spread, the properties of each agent, including health states s_{hea} and supply level s_{sup} will be updated. This section explains the three major steps of the simulator above in detail; the rest of our MPS model are explained in RL-related parts (for computing reward) and experiment (for government policy).

4.1 Distribution of Individuals to Facilities

Every day after agents make their decisions, they will be distributed into different facilities for calculation of contact.

Deciding the agent affiliation to facilities. We collect facility data from ArcGIS [ArcGIS, 2020], and we use synthesized US population dataset [Wheaton, 2012] for individual

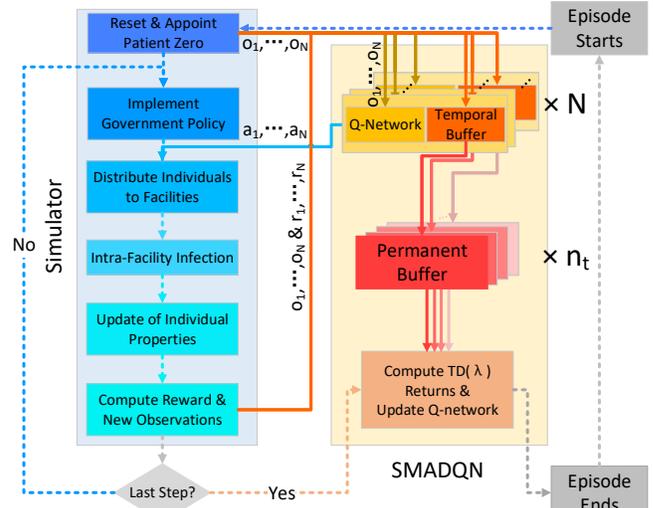


Figure 1: The flow chart of MPS and its interactions with SMADQN-controlled agents. The dotted lines are logical processes and solid lines are data flows.

agent data. Similar to the reality that individuals prefer to visit near facilities, we attach each available agent to the nearest facility of each type on the map within a distance upper bound. For simplicity, one individual is only attached to one building for each facility type. Another important thing we should mention is that age is crucial for modeling agents in epidemic outbreak since it is related with the type of facilities where one is active. In our model, the population is divided into four age groups, which are preschool children (0-9 years old), students (10-18), adults (19-64) and seniors (65+). We don’t consider preschool children and students’ after-school activities so they do not have a_{act} while adults and seniors have. Students are affiliated with schools and adults with workplaces while preschool children and seniors with neither.

Facility Attributions. How to distribute agents to facilities depends on agents’ actions and attributions of facilities. The two most important attributions about a facility are its capacity and type. For places other than hospitals, if there are more people going to a facility than its capacity, extra individuals will be uniformly randomly picked and kicked out from the facility that day. The type of facility decides the basic infection coefficient: I_F , which represents the relative infection probability in this type of facilities, and the minimum daily a_{act} for agents to join: $MinAct$. An agent will visit a facility with $MinAct = x$ when its $a_{act} \geq x$. For all non-compulsory facilities, $MinAct$ is assigned based on I_F : The higher I_F is, the higher the $MinAct$ is.

4.2 Intra-Facility Infection

The contact pattern of people inside a facility could be very complicated. In reality, some people may wander around while others may stay isolated. However, with the absence of GPS data, it would be too arbitrary to make any assumption on agents’ behavior patterns inside a facility. Even such assumptions can be made, the computational cost would be too high to keep track of agents’ trajectories. Therefore, we

simply assume that all infection happen in facilities and each person makes a constant number of contacts in a particular facility irrelevant of its capacity, and the contacts between agents are independently uniformly drawn. The probability p of an individual x infected in facility F with n people this day is: $p = \min(\beta I_F f_F p_x \sum_{y \in F} p_y / C_F, 1)$ where $p_x = g(a_{ge_x}, a_x)$ and $p_y = h(s_{y,hea}, a_y)$. In this formula, β is the overall hyper-parameter for infection rate. p_x is the factor related to the victim x , and p_y is the factor related to the contagious patient y . I_F is the basic coefficient of infection in facility F , f_F is the normal frequency of people going to facility F and C_F is the maximum capacity of facility F . p_x and p_y are calculated respectively by function g and h which are determined by the agent’s specific health state, action and age. For example, young and mask-worn people have less probability to get infected and asymptomatic, presymptomatic or mask-worn patients are less infectious.

There are two issues worth noting about this formula. First, it *smooths out* the infection probability by multiplying the frequency in normal life f_F onto the infectious probability. In another word, when an agent decides $a_{act} = x$, it will visit all facilities with $MinAct \leq x$ on that day being less infectious and susceptible to infection, and the reduction is inversely proportional to f_F , as if it will visit each facility with probability equal to $1/f_F$. Second, the formula above is an approximation to the real probability $p = \beta I_F f_F p_x (1 - \prod_y (1 - p_y)) / C_F$. This simplification accelerates the calculation process.

Detailed values of parameters, settings of functions and the algorithm deciding the agents’ affiliation are listed in the appendix. Moreover, community has a slightly different formula for infection, which is also stated in the appendix.

4.3 Update of Individual Properties

After the intra-facility disease spread, the two major properties of each agent i will be updated: supply level $s_{i,sup}$ and health state s_{hea} .

Supply Level. Supply level models the amount of stock one remains. It resembles a balance between sheltering-at-home and going out less and going shopping for essential goods to survive. Supply level drops monotonically from 1 to 0 with decreasing speed as people get more thrift with scarce supplies, and can only be reset to 1 when one of the agents in the household selects $a_{shop} \in \{\text{offline}, \text{online}\}$ and succeeds in shopping. Shopping offline is risky but will always succeed. Shopping online avoids infection, but the number of people it successfully serves every day is very limited, and failed online shopping brings nothing. Agents will receive increasing negative reward with decreasing supply level.

Health State. Each individual’s health state s_{hea} is updated at every simulation step, and the observation o_{hea} is decided by s_{hea} ; the transitions and correspondence are illustrated in fig. 2. Agents cannot directly access their health state s_{hea} , but can observe the observations o_{hea} instead.

We model the natural history of disease as a variant of the SEIR model [Dietz, 1976] with some changes: (1) presymptomatic, asymptomatic and immune states are added to better adapt to the real-world situation; (2) incubation and immune states adjacent to the symptomatic and asymptomatic state

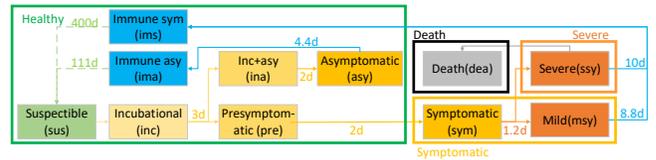


Figure 2: The natural history of disease based on the SEIR model. Each rectangle represents a health state (s_{hea}) and they are 4 boxes representing health observations (o_{hea}). The numbers on the edges are the expectation time of transition. The transition is also significantly infected by age; older people have much higher probability of being severe (from “sym” to “ssy”) or dead (from “ssy” to “dea”).

are separated to two states respectively for their infectivity and length differs; (3) the patients are divided into two sets: “mild” and “severe” to better model the hospitalization process. The mild patients do not need to be treated at hospital, and will recover automatically; the severe patients must be treated at hospital, otherwise it will die soon. Due to the lack of effective cure, treatment in hospital will not accelerate recovery of mild patients in our model.

5 SMADQN

In our settings, although sufficient data are generated after each episode, the value-policy iterations are very slow for simulating a whole episode costs long time. So, we choose a value-based algorithm like DQN as our base algorithm rather than other policy-based algorithms, whose policy only changes a little in the local area after each training step. But vanilla DQN is still problematic for such a massive MARL problem. First, for the property of partial observation, the consequence of actions may only be observed after several steps, so one-step TD learning is not sufficient. Second, it requires parameter sharing in policies to keep computation efficient with millions of agents. However, under deterministic policy, although the behaviour change for a single agent is incremental in training, the multi-agent system may experience dramatic change due to parameter sharing, which may lead to severe oscillation during training. Third, “Off-policy” methods may store some expired transitions in the replay buffer and disturb the training process while “on-policy” method may cause huge oscillation.

To alleviate these problems, we make changes below to form our algorithm. First, we use TD(λ) to train the Q-network. Like a standard DQN, we also have two Q-networks: the source network Q and the target network Q' , while Q is trained using MSE-loss to fit λ -returns. But to reduce computational cost, λ -returns are calculated using Q' only after an episode finished and before training Q , and Q' is only updated after training Q . So λ -returns do not have to be re-updated when Q is being trained. Second, like [Haarnoja et al., 2017], we also use Soft Q-learning where the probability of the individual i choosing a_i when observing o_i : $P(a_i|o_i)$ is directly proportional to $Q(a_i|o_i)/\alpha$, where α is the soft rate of Q-learning. Besides exploration provided by “Softness”, ϵ -greedy method is added to encourage extra exploration. And ϵ is set to decrease gradually. Third, we keep a balance between fully “off-policy” and fully “on-policy”: We build 2 types of replay buffers: the temporary one stores trajectories generated from the current episode, and the permanent one

stores trajectories used for DQN’s training. When an episode finished, β (in our implementation, $\beta = 1/3$) of trajectories in the permanent buffer are randomly chosen and replaced by random β of trajectories in the temporal buffer, and the other trajectories in the temporal buffer are thrown away.

In our epidemic model, individuals in the same age stage share parameters with each other. We use $n_t = 4$ sets of DQN network to model all individuals, one for each age group. More details of hyperparameters and pseudocodes are in the appendix.

6 Experiments

This section aims to validate that the MPS is detailed, authentic, explainable and adaptive to external changes, hence can be used for decision makers to derive public strategies upon. We first show that our model has the ability to present comprehensive microscopic details in section 6. Then, we show that our model can describe and predict the real-world dynamics by fitting the real world data on the spread of covid. Furthermore, in order to exploit our model and provide new candidate strategies for decision makers, we evaluate *information disclosure*, a novel mitigation strategy which gives agents information on the infection status of each facility except household with a 2-day delay to account for the time for information collection. Finally, we further explore our model by running *quarantine (QT)*, a widely-used government mitigation strategy, to show that our model is compatible with the main toolbox of real-world government and provide detailed information of the pandemic when executing this mitigation strategy for the decision makers.

Basic Settings of Dataset. We chose to evaluate and explore our model and SMADQN based on real-world data at Allegheny county in US, Pittsburgh. For simplicity, our model is closed in the sense that there is no incoming or outgoing agent (e.g., tourists). Based on ArcGIS [ArcGIS, 2020] and US population synthetic dataset [Wheaton, 2012], there are $N = 1188112$ residents (agents), among whom are 130451 preschool children (0-9 years old), 114867 students (10-18), 739183 adults (19-64) and 200011 seniors (65+). For the beginning of each episode, 2 agents are chosen to be symptomatic (health state is “sym”) and 8 agents to be exposed (health state is “inc”) near West Penn hospital, according to news [Doyle, 2020]. For all scenarios, the action of first 10 days is locked to be a fixed risky behaviour due to the latency of both the government and people; we chose 10 days because it is the interval between the first discovered cases in Allegheny and the day when stay-at-home order is put into effect (and this is why there is a sudden infectivity drop in most experiments). To generate a model that catches the local epidemic contact structure better and is prepared for fine-grained government strategy, we modeled 14 types of different facility and divided them into 4 categories by the amount of infection risk. Their distributions are illustrated in the last subfigure of section 6.

Calibration. The first and most important requirement for any model applied in real-world decision making is authenticity to ground truth. Thus, before testing different government strategies with our model, we first calibrated our model to fit the real-world data in Allegheny between March 14th, 2020

(when the first two cases in Allegheny are discovered) and the end of May, after which mass unpredictable protest took place. If not specified, all experiments in this section run for 80 days. And just as the real-world situation, the government only discloses the number of infected people on each day. There are three major criteria for our calibration: infected cases, hospitality and fatality. However, infected cases might be under-reported at the beginning of an outbreak due to insufficient testing. Therefore, we did not fit official reported cases, but infected cases inducted by hospitalized cases instead. The induction is based on data collected from US nationwide [Silva *et al.*, 2020] [NCIRD, 2020].

Most hyper-parameters about Covid-19 can be obtained by previous works on the disease. But for the lack of data about individuals’ behaviours, reward parameters are hard to summarize, and we could only hypothesize them. We left validation of reward parameters as future work. Among all reward parameters, R_{ill} is the most important factor for agent behavior control, which is the immediate penalty of agents fallen ill: Higher R_{ill} leads to more discreet agent behavior, such as wearing masks more often, stricter social distancing practice and less shopping frequency. Therefore, we mainly tuned R_{ill} in the calibration. We also sampled different R_{ill} in some experiments while fixing other reward parameters. Other hyper-parameters are left for sensitivity analysis in the appendix. Figure 3 show our calibration result. All curves in our paper are averaged from $3 \times$ last 10 episodes of 3 random seeds after training for 100 episodes if not stated otherwise. For total infected cases and total hospitalization cases, our model yields a result with the error between predicted summed infected cases and the actual one being 2.72% and the error for summed severe cases being 3.41%. The error is calculated and averaged only using 20th to 80th days because data were not accurate enough at the very beginning of the epidemic. The death cases are less coherent but still reasonable. This result proves that our model is authentic and can be explored for more government policies.

We modeled R_{ill} as $4500(\frac{80-d}{80})^4 + 21000(\frac{d}{80})^4$ where d is the day that first infection occurs. The reason is two-fold: first, it is clearly shown from the real-world news and data that people are (monotonically) increasingly aware of the virus, otherwise the number of infections will not drop; second, we tried to fit the data using R_{ill} functions with the relatively simple form of power function for better explainability and generalizability. Another thing worth noting is that there is a steep drop of infectivity from day 10, which resembles the effect of stay-at-home order. See appendix for details about the timeline of government strategy in our model.

Information Disclosure. Transparency is crucial for epi-

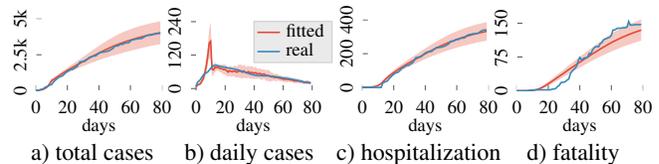


Figure 3: The result of calibration. The 4 sub-figures present total cases, daily cases, total hospitalization, and total deaths respectively. The real daily cases curve is smoothed for better readability.

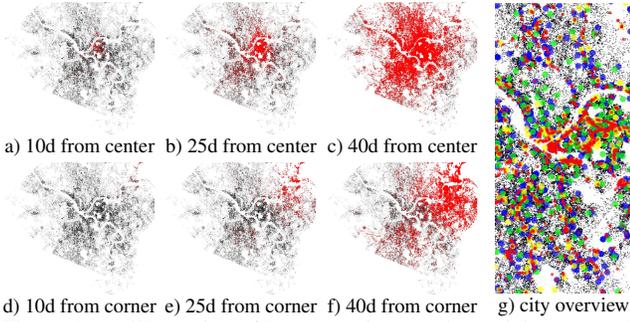


Figure 4: An illustration of how our microscopic model captures the local structure of the spread of disease and some other details. a) to c) are the global infection status of day 10, 25 and 40 starting from downtown, while d) to f) are the status of day 10, 25 and 40 starting from the edge of the map. The point is red if the household has at least 1 victim, and black otherwise. g) is an illustration of facility distribution in downtown area. The black points are households, and the red, yellow, green and blue points are the facilities with MinAct as 3, 2, 1 and 0.

demetic outbreak handling; not only will timely information about the infected cases alleviate panic, but ideally, when information is clear enough, agents will avoid infection sources for themselves and return to normal when the source disappears, thus minimizing the impact to economy and government costs. This proposes a new candidate strategy, *information disclosure*, to test the effect of pandemic control under ideal information coverage without any mandatory measures. In *information disclosure*, governments will keep monitoring new cases, and announce the probability of being infected in every non-household facility 2 days ago on a daily basis (we assumed that the government needs 2 days to collect information). Figure 5 shows the result of our experiments; the higher the R_{ill} is, which indicates that people are more serious about the disease, the lower the total number of cases is. The first row of fig. 5 clearly indicates that *information disclosure* is an effective strategy since it results in less cases than that without *information disclosure* in all four scenarios. It indeed reduces the number of daily cases, but the strategy may not be strong enough when individuals are less discreet since it still cannot flatten the curve when $R_{ill} \leq 3k$. The next rows illustrate how our RL model learns an adaptive response to different levels of risks. Higher observed risk leads to a drop in choosing $A_{act} = 3$ and offline shopping, as well as a rise in mask rate. Agents generally learned to avoid danger when the risk of infection increases, without explicitly setting a script for exact policy. More importantly, the learned policy has a higher average probability of taking risk with better control of epidemic, which illustrates that the behavior learned by RL is non-trivial in the sense that agents adopts a better policy to mitigate the pandemic while keeping a more normal life, which is possible because people in more dangerous areas are much more discreet. This effect slows down the spread dramatically.

Quarantine. In this scenario, besides information disclosure, the government will isolate every agent having been symptomatic for above 2 days (we assumed that it’s inevitable for governments to have delay). We tested two strategies with

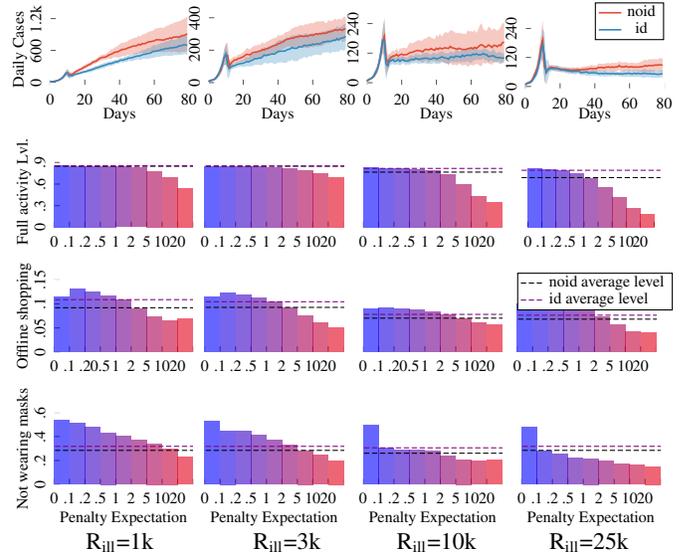


Figure 5: The experiment result of the *information disclosure* strategy. Every column is a different setting of R_{ill} .

different strength: The weaker strategy can only discover and quarantine a 3-day-infected people with probability = 1/3 each day. While The stronger one can discover it with probability = 1, and the strategy can also quarantine 40% agents that has directly infected by it as well, which needs contact tracing in reality. The isolated agent is excluded from the contact network. It cannot infect anybody upon isolation, including those in the same household, and will be released 9 days after recovery (in “ima” or “ims” state). Figure 6 shows the result of the three scenarios: no measure, weak quarantine, and strong quarantine. The first row shows that quarantine is very effective with most patients isolated from the contact network; even the weak quarantine can effectively control the pandemic (flatten the curve), and the strong quarantine eliminates pandemic within 80 days in half scenarios. The following two rows respectively shows the number of people isolated versus current number of cases under weak and strong quarantine strategy.

7 Conclusion

This paper introduced a microscopic epidemic simulator and the SMADQN algorithm to deal with millions of agents. We synthesized a comprehensive dataset for Allegheny county, US for evaluation, and explored two possible government policy candidates on our model: *information disclosure* and *quarantine*. Both proved useful for epidemic control.

A The Dataset based on Demographic Data of Pittsburgh

A.1 Data Source

Our dataset of Pittsburgh includes the following two types of information:

- Agents’ characteristics including their ages and trajectories/positions.

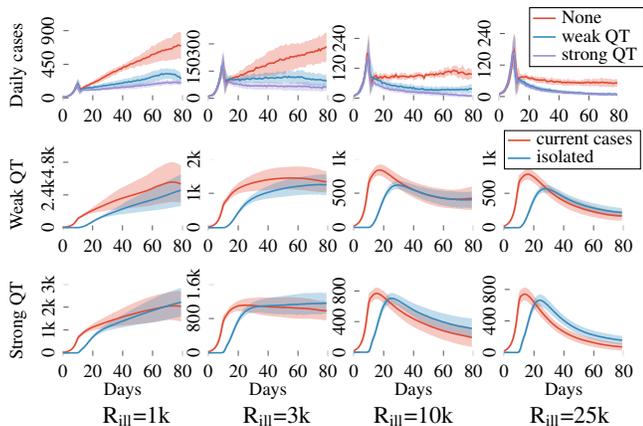


Figure 6: The experiment result of the *quarantine* (QT) strategy. Every column is a different setting of R_{ill} . The first row is the total cases of infection under three cases (red for no measure, blue for weak quarantine, and purple for strong quarantine). The second and third rows are the number of people isolated (blue) versus current number of cases (red) in weak QT and strong QT.

- Facility characteristics including locations, capacities and members of different facilities. For households, workplace and schools, we adopted the data from [Wheaton, 2012]. For recreational places, we utilized the data from [ArcGIS, 2020]. Regarding hospitals in Allegheny county, data from [of Health, 2019] was adopted.

A.2 Age Characteristics

Figure 7 shows the distribution of residents’ ages in our dataset. The distribution is synthesized based on the US population dataset [Bureau, 2020].

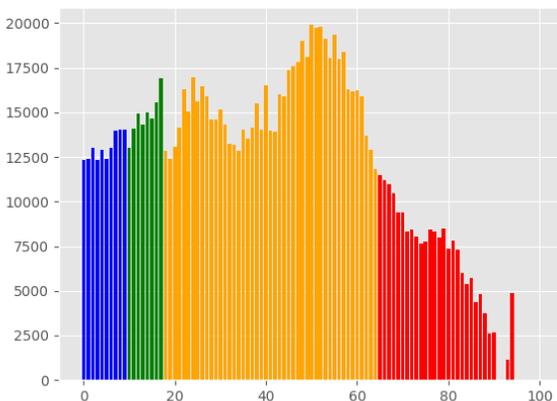


Figure 7: The distribution of residents’ ages in the US synthetic population dataset; the total population is 1, 188, 112. Blue, green, orange and red data stands for “chd” (children, 0-9 years old), “sch” (school students, 10-17 years old), “adu” (adults, 18-64 years old) and “rti” (retired people, 65+ years old) respectively.

A.3 Facility Characteristics

Table 1 lists the numbers of different types of facilities and their corresponding capacities (on average and maximally) in our synthesized dataset. In total, we have 14 types of facilities.

Facility	Number	Capacity (Average)	Capacity (Max)
Hospital	14	1577.43+384	8613+1542
Household	533919	2.22527	13
Workplace	38333	13.5461	7741
School	338	517.583	2239
Retail	601	247.977	1533
Supermarket	87	1900.97	5380
Community	358	3318.75	15889
Library	88	302.17	3069
Museum	78	77.8333	2347
Gym	193	110.782	1797
Restaurant	2691	201.458	3462
Stadium	3	3963.33	3976
Theatre	59	159.576	758
Cinema	36	367.917	3811

Table 1: The number, average capacity and maximum capacity of each type of facility (for hospitals, the capacity is in the “doctor + beds” format). We only model general hospitals in our dataset; specialized medical facilities such as rehabilitation centers and women’s hospitals are excluded.

A.4 Construction of the Contact Network

We followed the paradigm in [Aleta *et al.*, 2020] and [Kompella *et al.*, 2020] to construct the contact network among agents. The principles are: 1) all agents are affiliated with certain facilities; 2) at each simulation step, agents will be assigned into the affiliated facilities as long as they are not dead or hospitalized; 3) within the same facility, each pair of agents have equal probability to meet each other, i.e., equal probability to spread virus potentially; and 4) once an agent enters a facility, a minimum level of activity is required (see table 1 in the main manuscript for details about the requirements).

Following such principles, we first assigned each agent to a household based on the dataset in [Wheaton, 2012]. Similarly, agents above a certain age are also attached to schools and work places. However, the work in [Wheaton, 2012] didn’t provide other types of facilities as in our dataset (Table 1). To solve such a problem, we proposed to assign agents to facilities based on their distances to the facilities. The process is as follows:

- Step I: extract facility characteristics such as locations and capacities from [ArcGIS, 2020];
- Step II: determine the connection between agents and facilities based on their characteristics. For instance, children go to schools, adults go to workplaces, but the retired people go neither.
- Step III: establish the distribution of agents for each connected agent-facility group. We obtained such distributions by approximating the ratios of people visiting facilities in data from [Jones and Saad, 2019], [Hamm, 2020]

and [Statista, 2016]. Two major factors are considered here: capacities of the facilities and distances between agents and facilities. We prioritized agents closer to facilities, and facilities with larger capacities. For example, the priority radii for each facility type are 2km for conventional or retail stores, 5km for restaurants, gyms and supermarkets, 10km for theaters/cinemas, libraries and museums, and 35km for stadiums. Appendix A.4, appendix A.4, appendix A.4 and appendix A.4 illustrate agent allocation for different types of facilities. Note that hospitals are treated differently. We assume agents will go to the nearest hospital with remaining capacity > 0 considering the fact that life threaten is of highest priority.

B Parameters for the Agent Model

B.1 Hospitality and Fatality Rates

The death rate and hospitality rate used in our work are estimated based on the US population data [Bureau, 2020], and the CDC report on fatality cases and age-wise hospitalization rate per 100k people [for Disease Control and Prevention, 2020] (table 2 and table 3).

Age	proportion in infection	proportion in fatality
0-17	8.5%	0.06%
18-64	76.3%	20.71%
65+	15.2%	79.23%

Table 2: The proportion of each age group in overall infected and death cases[for Disease Control and Prevention, 2020] as of September 19th, 2020.

Age	Cumulative Hospitalization per 100k	Estimated US population in 2019
0-4	17.9	19756683
5-17	10.3	53462467
18-49	119.2	138216422
50-64	261.5	62925688
65+	472.3	54058263

Table 3: The estimated population, hospitalization rate and hospitalization cases for each age group.

B.2 Other Parameters

Table 5 shows the list of other parameters about the individual epidemiology.

C Details of Intra-Facility Infection

C.1 The calculation of p_x and p_y

As we mentioned in the main manuscript, $p_x = g(s_x, a_x)$ for a victim x and $p_y = h(s_y, a_y)$ for an infectious agent y are determined by multiple factors, such as whether the agent is wearing mask and the agent’s current health state. We considered the fact that people in different phases of the disease have different power of virus spreading.

More specifically, for all facilities except communities, $p_x = g(s_x, a_x)$ can be written as

$$p_x = p_{age,x} p_{mask,x},$$

and $p_y = h(s_y, a_y)$ as

$$p_y = p_{hs,y} p_{mask,y}.$$

In the formulae above, $p_{age,x}$ is the age factor of the victim x ; older people are more vulnerable to the infection and have bigger $p_{age,x}$. $p_{hs,y}$ is the factor for current health state of y ; people in pre-symptomatic phase and asymptomatic patients have limited power to spread virus. For any agent z , $p_{mask,z} = 0.4$ if the agent is wearing a mask ($a_{z,mask} = \text{mask}$), and 1 otherwise. For example, if both x and y are wearing masks, then the probability of infection is reduced to 16% compared to that without masks on both sides. For $p_{age,x}$ and $p_{hs,y}$, we have

$$p_{age,x} = \begin{cases} 0.4 & x \text{ is "chd"} \\ 0.38 & x \text{ is "sch"} \\ 0.8175 & x \text{ is "adu"} \\ 0.81 & x \text{ is "rtr"} \end{cases}$$

and

$$p_{hs,y} = \begin{cases} 0.12 & \text{Health state is pre-symptomatic} \\ 0.31 & \text{Health state is asymptomatic} \\ 1 & \text{Health state is in \{"sym", "msy" and "ssy"\}} \\ 0 & \text{otherwise} \end{cases}$$

As for other facility-related parameters for intra-facility infection, table 6 shows the frequency and basic coefficient of each type of facilities as well as their $MinAct$.

C.2 Special Rules for Community in Intra-Facility Infection

In our model, community infection stands for the probability of getting infected by walking by or chatting with the infected agents in open space. Infections on public transportation and other places are not considered here. The reason for us to distinguish community infection from other facilities is that the risk of community infection and agents’ activity levels are co-related. It is not accurate enough to model the risk of community infection by a single threshold as in other facilities.

Therefore, we let the activity level A_{act} as a factor of the spreading power in community infection. For example, if the infected individual chooses $A_{act} = 2$ and the victim chooses $A_{act} = 3$, then the probability of infection calculated from the formula in the paper should multiply $\frac{2}{2} * \frac{3}{2} = 1.5$. A cautious person will avoid community infection if and only if he/she chooses $A_{act} = 0$, which stands for staying at home as much as possible.

C.3 The Calculation of Contact Trace

Strong quarantine in our experiment requires contact tracing, namely, if a case is discovered by the government, “other people he/she has physically contacted” must be quarantined as

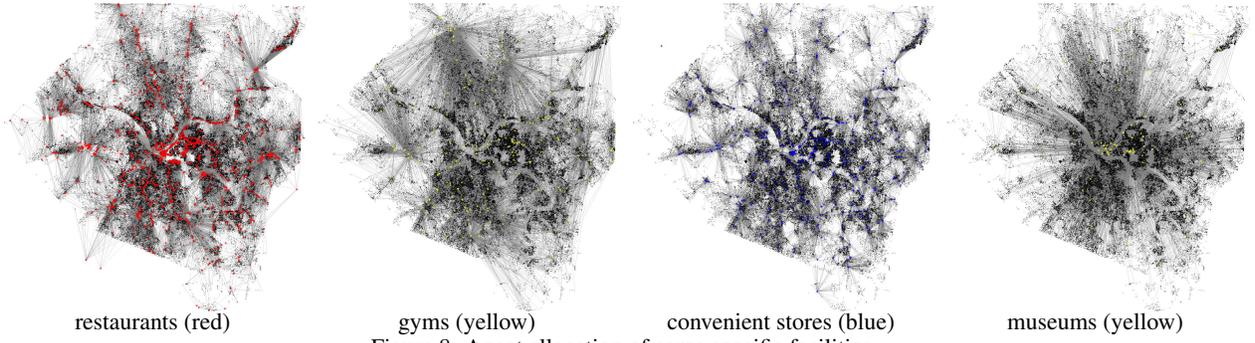


Figure 8: Agent allocation of some specific facilities.

Age	Estimated Infection	Estimated Cumulative Hospitalization Cases	Estimated Fatality	Severity Rate	Unrevised Death Rate
0-17	574796	9043	120	1.57%	0.021%
18-64	5159636	329304	41256	6.38%	0.8%
65+	1027869	255317	157831	24.84%	15.36%

Table 4: The calculated death rate and hospitalization rate for each age group in the US by Sept. 19th, 2020. We timed 1.25 to the final fatality rate for there exists on-going cases that will be dead in the future.

well. To calculate the contact trace in the intra-facility infection process, we assume the patient-victim pair (*i.e.* contact trace) is distributed proportionally to the virus-spreading power of all patients in the same facility every day. For example, if two patients *A* and *B* infect 3 victims in a workplace on a particular day, and *A*'s virus-spreading power is twice of *B*'s, then 2 people (randomly drawn) will be modeled as infected by *A*.

D Reinforcement Learning Agent Setup

D.1 detailed settings of the MA-POMDP

Observations: in experiments without information disclosure, an agent *i*'s observation is $o_{i, \text{noinfo}} = o_{i, \text{hea}} \times o_{i, \text{rel}} \times o_{i, \text{sup}} \times o_{\text{city}}$. $o_{i, \text{hea}}$ indicates the agent's health states, which is a one-hot vector with a dimension of 5¹. In our experiment, for simplicity, we assumed that all people without medical checks will believe themselves as infected by COVID-19 only when they have symptoms, and before that, people do not take medical tests. $o_{i, \text{rel}}$ indicates agent *i*'s relatives' health states (a relative is defined as another person in the same household), which is a real number and equals to the probability of agent *i* being infected by the relatives. $o_{i, \text{sup}} \in [0, 1]$ represents agent *i*'s supply level². Finally, o_{city} is the number of infected cases in the city, normalized by dividing 1000.

In experiments with information disclosure³ the observation of agent *i* is $o_{i, \text{info}} = o_{i, \text{noinfo}} \times o_{i, \text{sur}}$. We separate all

¹Note that although there are 11 types of health states, the agent cannot distinguish them all before medical check. For a comprehensive relationship between $o_{i, \text{hea}}$ and the health states, one can refer to figure 2 in our main manuscript.

²If the supply is not replenished, we let the supply level $o_{i, \text{sup}}$ drop at a decreasing rate. More specifically, denote the supply level as *L*, then $L = \max(0, 1 - (\frac{d}{21})^2)$, where *d* is the number of days since last replenishment of supply.

³We assumed that the government can only know and disclose

types of facilities into 4 groups and each group of facilities has the same level $\text{MinAct} \in \{0, 1, 2, 3\}$. The component of each group can be seen in table 6 in the appendix. $o_{i, \text{sur}} \in \mathbb{R}^4$ indicates severity of infections in all four groups of facilities that agent *i* may visit. We assumed that severity approximately equals to the probability of agent *i* being infected in those facilities.

All observations are concatenated together and feed into the Q-network. We use different sets of hyper-parameters to learn different policies.

Actions: Agent *i*'s action space is discrete, represented by $a_i = a_{i, \text{mask}} \times a_{i, \text{act}} \times a_{i, \text{shop}}$. $a_{i, \text{mask}} = \{\text{mask}, \text{no_mask}\}$ indicates wearing a mask or not. $a_{i, \text{act}} = \{0, 1, 2, 3\}$ indicates the activity level for other public facilities except retail stores. $a_{i, \text{shop}} = \{\text{no_shopping}, \text{shopping_online}, \text{shopping_offline}\}$ indicates the ways for shopping to replenish an agent's supply. In each simulation step (a day), there will be only 17,000 people, or roughly $\frac{1}{70}$ of the total population randomly chosen from the pool of agents with $a_{i, \text{act}} = \text{"shopping online"}$. The number $\frac{1}{70}$ is estimated from [Bishop, 2020].

Rewards: An agent rewards include rewards on activity levels, rewards on health states, rewards for wearing masks, rewards for offline shopping, and rewards related to the supply level.

We let the activity rewards $R_{\text{act}}(a_i)$ be positively proportional to the activity level if no symptom shows, *i.e.*, $R_{\text{act}}(a_i) = \alpha_{\text{act}} a_i$. Once infected, an one-time negatively high reward R_{ill} will be posted on the agent. $R_{\text{mask}} = r_{\text{mask}} < 0$ is assigned for wearing a mask, $R_{\text{shop}} = -1$ is assigned for selecting offline shopping, and $R_{\text{eth}} = -(\alpha_{\text{act}} a_i + r_{\text{mask}})$ (the ethical penalty) is given for high activity level or not wearing mask with symptoms. Also, an agent gets penalty for low supply level. More specifically, it gets a negative re-

information of people infected 2 days ago.

Notation	Meaning	Value	Source
p_sev2cri_nhos	The daily probability of death for severe patients out of hospital(0-18 yrs old/18-65/65+)	0.6,0.8,1	estimated
p_sev2rec	The daily probability of recovery for severe patients in hospital	1/10	[Aleta <i>et al.</i> , 2020]
p_inc2pre	The daily probability of developing from "inc" state to "pre" or "ina" state	1/3	[Aleta <i>et al.</i> , 2020]
p_rec_sym	The daily probability of recovery from mild symptom	1/8.8	[Tambri Housen and Sheel, 2020]
p_hos	The daily probability of developing from "sym" state to "msy"/"ssy" state	1/1.2	[Lauer <i>et al.</i> , 2020]
p_deimm_asy	The daily probability of losing immunity after being asymptomatic	0.009	[Long <i>et al.</i> , 2020]
p_deimm_sym	The daily probability of losing immunity after being symptomatic	0.0025	[Long <i>et al.</i> , 2020]
asy_infect_rate	The infectivity of asymptomatic patients (symptomatic is 1)	0.31	[Aleta <i>et al.</i> , 2020]
pre_infect_rate	The infectivity of pre-symptomatic patients	0.12	[Aleta <i>et al.</i> , 2020]
asy_pop_rate	The proportion of patients that turn out to be asymptomatic.	0.25	calibrated

Table 5: The individual epidemiology parameter list. All daily probabilities are geometrical distribution and independent for each day.

Facility	I_F	f_F	MinAct
Hospital	0	N/A	0
Household	0.23	1	0
Workplace	0.14	5/7	0
School	0.21	5/7	0
Retail	0.09	1	0
Supermarket	0.09	1	0
Community	0.0075	1	1
Library	0.12	10.5/365	2
Museum	0.12	2.5/365 * 1/0.54	2
Gym	0.15	0.47	2
Restaurant	0.21	4.2 / 7	3
Stadium	0.42	4.7/365 * 1/0.17	3
Theatre	0.42	3.8/365 * 1/0.35	3
Cinema	0.42	5.3/365 * 1/0.59	3

Table 6: The basic coefficient(I_F), normal frequency(f_F) and minimal a_{act} (MinAct) for joining the calculation of each facility(Fac). Some data (fractions) are corrected by the proportion of active person in our model of Allegheny. Data are based on [Aleta *et al.*, 2020][Hamm, 2020][Jones and Saad, 2019][Statista, 2016]

ward of $r_{sup} = -\frac{1}{0.58}(1 - L)$ with L as the supply level. The constant $\frac{1}{0.58}$ is calculated to make a rational agent's shopping frequency at around 7–8 days, beyond which the incentive to avoid infection in offline shopping by R_{shop} is overwhelmed.

In our settings, we assumed $\alpha_{act} = 1$ and $r_{mask} = 0.1$. We tried different values for R_{ill} to generate different settings where people care about their health in different degrees. We left finding a realistic set of more hyper-parameters such as α_{act} and r_{mask} via inverse RL as future work.

Note that by the one-time penalty R_{ill} , the RL training in our work is naturally smoothing. The smoothing process is two-fold: 1) one-time penalty is smoothed into each simulation step; for example, if a healthy person goes to places where the probability of infection is 10%, then he would receive $0.1R_{ill}$; 2) people don't go to every facility at each step, but we smoothed it by correcting the factor with frequency. For example, if people goes to workplaces 5 times per week and libraries 10 times per year, we will assume that they go to such places every day, with the virus-spreading power multiplied by $\frac{5}{7}$ and $\frac{11}{365}$.

D.2 details of SMADQN

Algorithm 1 and table 7 are, respectively, the pseudocode and hyper-parameters of the SMADQN algorithm.

Hyper-parameters	Value
minimum of ϵ (ϵ_m)	0.9
maximum of ϵ (ϵ_M)	1.2
step-size of ϵ (ϵ_s)	0.1
Soft rate (α)	1 / 3
training optimizer	Adam
Learning rate	$0.01x + 0.001(1 - x)$, $x = n_{episode}/20 - 1$ clipped to $[0, 1]$
# mini-batch (n_b)	100
Adam step-size	1e-5
Discount rate (γ)	0.9
GAE parameter (λ)	0.9
episode length (T)	80

Table 7: The hyper-parameters in the SMADQN algorithm

E Major Covid-Related Events of Allegheny County and Corresponding Government Policies

To better simulate the real-world infection, we collected major news related to the spread of Covid-19 in Allegheny County, US. Table 8 listed the major events since Covid-19 began to spread in Allegheny. The collected events helped us calibrate the hyper-parameters of our model from two major aspects: 1) we fixed all actions of individuals and implemented no government policy for the first 10 days in all experiments to simulate the delay of the government and regular people; 2) we chose 80 days as the length of one episode to match the fact that Covid-19 fighting in Allegheny achieved a stage of success in the first 80 days since the first case was discovered. Table 8 shows the major events and the corresponding government policies in our model.

Algorithm 1 SMADQN

```
1: Randomly initialize source Q network for each agent type
    $t$ :  $Q_t(o, a|\theta_t)$  with weights  $\theta_t$ .
2: Initialize target Q network for each agent type  $t$ :  $Q'_t$  with
   weights  $\mu_t \leftarrow \theta_t$ 
3: Initialize permanent buffer for each agent type  $t$ :
    $R_{perm,t}$ .
4: Initialize a standard normal distribution generator  $G$ 
5: Initialize the greedy threshold  $\epsilon = \epsilon_m$ 
6: for episode  $i = 1, 2, \dots, +\infty$  do
7:   Initialize temporal buffer for each agent type  $t$ :
      $R_{temp,t}$ 
8:   Receive initial observation for each individual  $i$  of
     each type  $t$  at step 1:  $o_{t,i,1}$ 
9:   for  $j = 1, \dots, T$  do
10:    for each individual  $i$  of each type  $t$  do
11:      Select action of this individual according
      to its source Q network:  $a_{t,i,j} \propto$ 
       $Q_{t,i}(a_{t,i,j}, o_{t,i,j}|\theta_t)/\alpha$ 
12:      use  $G$  to generate a random number  $g$ 
13:      if  $g \geq \epsilon$  then
14:        Resample  $a_{t,i,j}$  uniformly
15:      end if
16:    end for
17:    Execute actions and observe reward  $r_{t,i,j}$  nad ob-
    serve new observation  $o_{t,i,j+1}$ 
18:  end for
19:  for each type  $t$  do
20:    if  $episode = 1$  then
21:      Store all individuals' trajectories:
       $(o_{t,i,1}, a_{t,i,1}, r_{t,i,1}, o_{t,i,2}, a_{t,i,2}, r_{t,i,2}, \dots, o_{t,i,T},$ 
       $a_{t,i,T}, r_{t,i,T})$  in  $R_{temp,t}$ 
22:    else
23:      Replace  $\beta$  of trajectories in  $R_{perm,t}$  with trajec-
      tories in  $R_{temp,t}$ 
24:    end if
25:    Update  $Q'_t$ :  $\mu_t \leftarrow \theta_t$ 
26:    Use  $Q'_t$  to update GAE values for each  $o_{t,i,j}$  in
     $R_{perm,t}$  as  $y_{t,i,j}$ 
27:    Separate data in  $R_{perm,t}$  to  $n_b$  mini-batches, and
    train  $Q_t$  with MSE using  $y$  as target.
28:  end for
29:   $\epsilon \leftarrow \min(\epsilon + \epsilon_s, \epsilon_M)$ 
30: end for
```

F Computational Feasibility

Our experiment code was developed with multi-threaded C++ for simulation and Python for MARL training and data processing. We connected the two parts with Cython. We run our experiments on a server with a CPU: Intel(R) Core(TM) i9-9940X CPU @ 3.30GHz (14 cores) and a GPU: RTX 2080 Ti. A typical step usually took around 15-20 seconds in our experiments, and about half an hour for a complete epoch. The whole program needs 40G RAM and 3G VRAM. The whole algorithm has $O(n)$ time complexity and $O(n)$ space complexity, where n is the number of agents. We use multi-thread to reduce the constant of $O(n)$. Most experiments were

Day	Event	Government Policy (Capacity Restraint)
0	First two cases are discovered in Allegheny	N/A
10	The stay-at-home order is in effect	workplace 25%, supermarket, community & retail 100%, others 0% (capacity)
62	Allegheny move to "yellow" phase for reopening	workplace 50% community 100% supermarket & retail 100% restaurant 25%
80-100	Massive protest	N/A
82	Allegheny move to "green" phase for reopening	workplace 75% community 100% supermarket & retail 100% school 0% others 50%
110	Temporal ban on restaurants	workplace 75% community 100% supermarket & retail 100% restaurant 0% others 50%
116	Soften regulations for restaurants	workplace 75% community 100% supermarket & retail 100% restaurant 10% others 50%
123	Soften regulations for restaurants	workplace 75% community 100% supermarket & retail 100% restaurant 35% others 50%

Table 8: Major events related to Covid-19 outbreak in Allegheny County [All, 2020] and the corresponding government policies. Day 0 is March 14th, 2020. Although we did not simulate the situation after 80 days, we still assigned government policies for them.

trained for 100 episodes to guarantee convergence, which usually has much redundancy in practice. Hence, the training process of our experiments could be even much faster with less episodes.

G Boosting Training with Policy Transfer

A useful property of SMADQN is that it does not need to be trained from scratch for each experiment. Instead, we show in this section that the policy can be transferred between different experiments to boost training: 1) train a policy on the settings without any government control, and 2) transfer the policy on a setting with much strong control and fine-tune it for a few episodes.

Figure 9 illustrates the performance of SMADQN with different settings under strong government control, i.e., the real-life government strategy used in calibration. Via the results, we can see that the policy transfer worked well. We use the similarity of average daily cases as the performance metric.

We can see that with policy transfer ("None+finetune"), SMADQN quickly yielded a policy which behaved very sim-

ilarly to the policy trained under strong control from scratch for 100 episodes (“Strong”). The outcome is more similar than the policy trained from scratch for 20 episodes (“Strong 20 episode”). Such observation indicates that policy transfer with SMADQN is more efficient than training from scratch. Direct application of policy trained on non-control settings yielded a very different outcome (“None”), which shows the necessity of policy transfer. Moreover, the difference is coherent with the common knowledge that being less serious about the virus will lead to much poorer control effect of the pandemic.

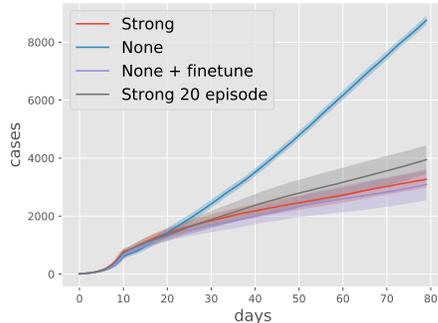


Figure 9: The averaged number of summed cases with strong government control (the same with the calibration). “Strong” stands for training with strong control for 100 episodes”; “None” stands for training with no governmental control for 100 episodes; “None+finetune” stands for training with no control for 80 episodes and with strong control for 20 episodes; “Strong 20 episode” stands for training with strong control for 20 episodes from scratch.

H Sensitivity Analysis

We show that our MPS model and the SMADQN algorithm behave reasonably in the sensitivity analysis by varying two hyper-parameters: the infection rate β and the penalty of wearing mask R_{mask} . The results are shown in fig. 10 and fig. 11. All data in the sensitivity analysis was the average of last 10 episodes.

Figure 10 shows the daily number of cases under different settings of β . $\beta_0 = 15.8$ is the default β in our experiments. All RL agents were trained without government control strategies. We can see that the daily cases increased monotonically with β , which shows that the MPS is stable and can produce reasonable results for different β . Furthermore, for all three β s, the epidemic was controlled to some extent instead of growing exponentially, which means that our SMADQN can handle a wide range of hyper-parameters related to the real-world environment.

Figure 11 depicts the daily number of cases under different settings of R_{mask} ($R_{mask} = -0.1$ in a default experiment setting). RL agents were all trained without government control strategies. The daily cases increased monotonically with more harsh R_{mask} ’s penalties, which shows that our SMADQN can handle a wide range of reward parameters and yield reasonable policies in coherent with common knowledge. Hence, it indicates that our reward parameters can successfully capture the influence of individual values towards the macroscopic development of pandemic. Moreover,

the SMADQN algorithm can also be utilized as a feasible tool for real-world decision makers.

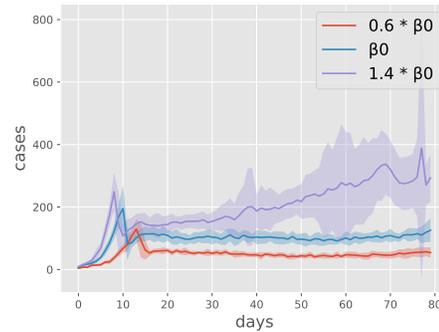


Figure 10: Daily cases with different β and no government control strategy.

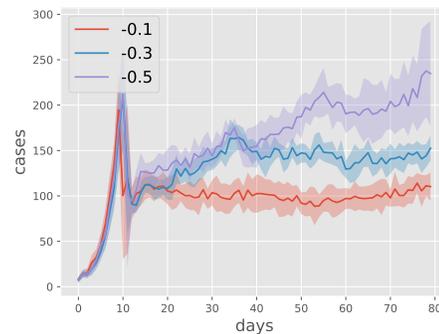


Figure 11: Daily cases with different penalty $R_{mask} \in \{-0.1, -0.3, -0.5\}$ for wearing masks.

I Code and Data Availability

All codes and data are accessible in <https://github.com/recordmp3/Microscopic-epidemic-model/settings>, and the format of the data can be seen in the code.

References

- [Aleta *et al.*, 2020] Alberto Aleta, David Martın-Corral, Ana Pastore y Piontti, Marco Ajelli, Maria Litvinova, Matteo Chinazzi, Natalie E. Dean, M. Elizabeth Halloran, Ira M. Longini Jr, Stefano Merler, Alex Pentland, Alessandro Vespignani, Esteban Moro, and Yamir Moreno. Modeling the impact of social distancing, testing, contact tracing and household quarantine on second-wave scenarios of the covid-19 epidemic. *Nature Human Behaviour*, 4:964–971, 2020.
- [All, 2020] Stay updated - allegheny county, 2020. [<https://www.alleghenycounty.us/Health-Department/Resources/COVID-19/Stay-Updated.aspx>; accessed 27-October-2020].
- [ArcGIS, 2020] ArcGIS. Arcgis business analyst, 2020. [<https://bao.arcgis.com/esriBAO/index.html>; accessed 25-August-2020].
- [Bishop, 2020] David Bishop. June 2020 online grocery scorecard: Growth in sales & hh penetration

- continues, 2020. [<https://www.brickmeetsclick.com/june-2020-online-grocery-scorecard-growth-in-sales-hh-penetration-continues>; accessed August-15-2020].
- [Bureau, 2020] US Census Bureau. National population by characteristics: 2010-2019, 2020. [<https://www.census.gov/data/tables/time-series/demo/popest/2010s-national-detail.html>; accessed 15-October-2020].
- [Chang *et al.*, 2020] Sheryl L. Chang, Nathan Harding, Cameron Zachreson, Oliver M. Cliff, and Mikhail Prokopenko. Modelling transmission and control of the covid-19 pandemic in australia. *Nature Communications*, 11(5710), 2020.
- [Dietz, 1976] K. Dietz. The incidence of infectious diseases under the influence of seasonal fluctuations. In Jürgen Berger, Wolfgang J. Bühler, Rudolf Repges, and Petre Tautu, editors, *Mathematical Models in Medicine*, pages 1–15, Berlin, Heidelberg, 1976. Springer Berlin Heidelberg.
- [Doyle, 2020] Patrick Doyle. Allegheny county announces first two cases of covid-19, 2020. [<https://www.witf.org/2020/03/14/allegheny-county-announces-first-two-cases-of-covid-19/>; accessed 15-October-2020].
- [Eilersen and Sneppen, 2020] Andreas Eilersen and Kim Sneppen. Cost–benefit of limited isolation and testing in covid-19 mitigation. *Nature Scientific Reports*, 10(18543), 2020.
- [for Disease Control and Prevention, 2020] Centers for Disease Control and Prevention. Cdc covid data tracker, 2020. [<https://covid.cdc.gov/covid-data-tracker/#demographics>; accessed 27-October-2020].
- [Haarnoja *et al.*, 2017] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *ICML*, 2017.
- [Hamm, 2020] Trent Hamm. Don’t eat out as often, 2020. [<https://www.thesimpledollar.com/save-money/dont-eat-out-as-often>; accessed 15-October-2020].
- [Hoertel *et al.*, 2020] Nicolas Hoertel, Martin Blachier, Carlos Blanco, Mark Olfson, Marc Massetti, Marina Sánchez Rico, Frédéric Limosin, and Henri Leleu. A stochastic agent-based model of the sars-cov-2 epidemic in france. *Nature Medicine*, page 1417–1421, 2020.
- [Jones and Saad, 2019] Jeff Jones and Lydia Saad. Gallup news service december wave one final topline, 2019. [Timberline: 937008; JT: 335; Princeton Job number: 19-12-021; accessed 13-October-2020].
- [Kompella *et al.*, 2020] Varun Kompella, Roberto Capobianco, Stacy Jong, Jonathan Browne, Spencer Fox, Lauren Meyers, Peter Wurman, and Peter Stone. Reinforcement learning for optimization of covid-19 mitigation policies. *arXiv preprint arXiv:2010.10560*, 2020.
- [Lauer *et al.*, 2020] Stephen A. Lauer, Kyra H. Grantz, Qifang Bi, Forrest K. Jones, Qulu Zheng, Hannah R. Meredith, Andrew S. Azman, Nicholas G. Reich, and Justin Lessler. The incubation period of coronavirus disease 2019 (covid-19) from publicly reported confirmed cases: Estimation and application. *Annals of Internal Medicine*, 2020.
- [Liu, 2020] Changliu Liu. A microscopic epidemic model and pandemic prediction using multi-agent reinforcement learning. *ArXiv*, abs/2004.12959, 2020.
- [Long *et al.*, 2020] Quan-Xin Long, Xiao-Jun Tang, Qiu-Lin Shi, Qin Li, Hai-Jun Deng, Jun Yuan, Jie-Li Hu, Wei Xu, Yong Zhang, Fa-Jin Lv, Kun Su, Fan Zhang, Jiang Gong, Bo Wu, Xia-Mao Liu, Jin-Jing Li, Jing-Fu Qiu, Juan Chen, and Ai-Long Huang. Clinical and immunological assessment of asymptomatic sars-cov-2 infections. *Nature Medicine*, 26:1200–1204, 2020.
- [Meloni *et al.*, 2011] Sandro Meloni, Nicola Perra, Alex Arenas, Sergio Gómez, Yamir Moreno, and Alessandro Vespignani. Modeling human mobility responses to the large-scale spreading of infectious diseases. *Scientific reports*, 1:62, 2011.
- [NCIRD, 2020] NCIRD, 2020. [National Center for Immunization and Respiratory Diseases (NCIRD), Division of Viral Diseases, <https://www.cdc.gov/coronavirus/2019-ncov/covid-data/investigations-discovery/hospitalization-death-by-age.html>; accessed 15-October-2020].
- [of Health, 2019] Pennsylvania Department of Health. Hospital reports, 2019. [<https://www.health.pa.gov/topics/HealthStatistics/HealthFacilities/HospitalReports/Pages/hospital-reports.aspx>; retrieved 2020-08-31].
- [Qu and Li, 2019] Guannan Qu and Na Li. Exploiting fast decaying and locality in multi-agent mdp with tree dependence structure. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 6479–6486. IEEE, 2019.
- [Qu *et al.*, 2020] Guannan Qu, Adam Wierman, and Na Li. Scalable reinforcement learning of localized policies for multi-agent networked systems. In *Learning for Dynamics and Control*, pages 256–266. PMLR, 2020.
- [Silva *et al.*, 2020] Petrônio CL Silva, Paulo VC Batista, Hélder S Lima, Marcos A Alves, Frederico G Guimarães, and Rodrigo CP Silva. Covid-abs: An agent-based model of covid-19 epidemic to simulate health and economic effects of social distancing interventions. *arXiv preprint arXiv:2006.10532*, 2020.
- [Statista, 2016] Research Department Statista. How often people work out at the gym in us 2016, 2016. [<https://www.statista.com/statistics/638978>; accessed 15-October-2020].
- [Tambri Housen and Sheel, 2020] Amy Elizabeth Parry Tambri Housen and Meru Sheel. How long are you infectious when you have coronavirus?, 2020. [<https://theconversation.com/how-long-are-you-infectious-when-you-have-coronavirus-135295>; accessed 15-October-2020].
- [Wheaton, 2012] WD Wheaton. 2009 us synthetic population ver. 2, 2012.

[Yang *et al.*, 2018] Yaodong Yang, Lantao Yu, Yiwei Bai, Ying Wen, Weinan Zhang, and Jun Wang. A study of ai population dynamics with million-agent reinforcement learning. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages 2133–2135, 2018.

[Zheng *et al.*, 2017] Lianmin Zheng, Jiacheng Yang, Han Cai, Weinan Zhang, Jun Wang, and Yong Yu. Magent: A many-agent reinforcement learning platform for artificial collective intelligence. *arXiv preprint arXiv:1712.00600*, 2017.