

Location-Based Social Simulation for Prescriptive Analytics of Disease Spread

Joon-Seok Kim¹, Hamdi Kavak², Chris Ovi Rouly³, Hyunjee Jin¹,
Andrew Crooks^{1,2}, Dieter Pfoser¹, Carola Wenk³, Andreas Züfle¹

¹Department of Geography and Geoinformation Science, George Mason University, USA

²Department of Computational and Data Sciences, George Mason University, USA

³Department of Computer Science, Tulane University, USA

{jkim258,hkavak,hjin6,acrooks2,dpfoser,azufle}@gmu.edu,{orouly,cwenk}@tulane.edu

Abstract

Human mobility and social networks have received considerable attention from researchers in recent years. What has been sorely missing is a comprehensive data set that not only addresses geometric movement patterns derived from trajectories, but also provides social networks and causal links as to why movement happens in the first place. To some extent, this challenge is addressed by studying location-based social networks (LBSNs). However, the scope of real-world LBSN data sets is constrained by privacy concerns, a lack of authoritative ground-truth, their sparsity, and small size. To overcome these issues we have infused a novel geographically explicit agent-based simulation framework to simulate human behavior and to create synthetic but realistic LBSN data based on human patterns-of-life (i.e., a geo-social simulation). Such data not only captures the location of users over time, but also their motivation, and interactions via temporal social networks. We have open sourced our framework and released a set of large data sets for the SIGSPATIAL community. In order to showcase the versatility of our simulation framework, we added disease a model that simulates an outbreak and allows us to test different policy measures such as implementing mandatory mask use and various social distancing measures. The produced data sets are massive and allow us to capture 100% of the (simulated) population over time without any data uncertainty, privacy-related concerns, or incompleteness. It allows researchers to see the (simulated) world through the lens of an omniscient entity having perfect data.

1 Introduction

The detailed investigation of collective social phenomena requires an inordinate amount of data given that the patterns underlying human behavior can be quite complex and are as such hard to predict. This is summed up quite aptly by the late Nobel laureate Murray Gell-Mann: “*Think how hard physics would be if particles could think.*” To some extent, research on location-based social networks (LBSNs) attempts to grapple this challenge by leveraging social network data for predictive tasks such as Point-of-Interest recommendation [25, 27, 28], social link prediction [21], and location prediction [6]. The major challenge, however, is that comprehensive real-world LBSN data sets are hardly available due to privacy implications. Also, because such data are considered operational data, businesses are unwilling to make them publicly available or share them. The largest publicly available LBSN data set is the Gowalla data set [6] having 36M check-ins. But, after removing users with less than 15 check-ins, and removing locations with less than 10 visitors, from Gowalla, only 18.7k users and 1.29M

check-ins remain [18]. Distributed over 20 months, that is an average of only 2.1k checkins per day. Distributed globally this leaves only a hand-full of check-ins per day per city, hardly enough to model, explain and predict mobility. A recent study by Li et al. [17] concludes that “*Researchers working with LBSN data sets are often confronted by themselves or others with doubts regarding the quality or the potential of their data sets*” and that “*it is reasonable to be skeptical*”[17].

Towards addressing this challenge, we have developed an agent-based simulation framework that exhibits realistic social behavior based on the data about real-world phenomena and social science theories. We simulate plausible numbers of agents over years of simulation time and potentially generate LBSN data for entire generations. The simulation generates high-fidelity LBSN data sets containing complete location and temporal social network data without uncertainty collected over long periods of time. We have published this agent-based simulation framework, ran many simulations, and made available many such data sets [13]. LBSN is certainly important but one of the many potential uses of our data sets. Given the current pandemic crisis, there is a potential to use our framework for disease spread simulations.

From cholera outbreaks that continue to plague certain parts of the world to the more recent Ebola, Zika and now the COVID-19 pandemic, data has been crucial with respect to monitoring and understanding such diseases. John Snow [22], for example, traced the source of a cholera epidemic to a specific water pump by talking to people infected with cholera. His analysis not only led to the foundation of modern epidemiology [19], but served as an exemplar of why the geospatial data aspect is important when it comes to exploring the spread of a disease. Recent disease outbreaks have benefited from advances in technology (e.g., mobile phone records and ubiquitous internet) and contact tracing to understand how diseases spread [3]. As publicly available data is not always sufficient to the understanding of such a complex phenomena, the work presented in this article aims at augmenting data analysis with modeling to better comprehend and predict the evolution and impact of diseases such as COVID-19. Agent-based models are perfect fit for this purpose, compared to mathematical models [5, 7, 24], as populations are heterogeneous and their infection risk is vastly different according to their mobility, age, pre-existing conditions etc. Within this style of model, individuals are specifically accounted for and focus on the dynamic interactions between individual and their environment are captured, such as the daily patterns-of-life (PoL) (e.g., [9, 13]).

The purpose of this newsletter article is to: (1) introduce the large LBSN data sets mentioned above to the SIGSPATIAL community, (2) showcase the potential for simulating disease outbreaks, and (3) study the effects of a range of mitigation measures on disease spread. In the remainder of this article we introduce our agent-based simulation framework that can simulate daily patterns-of-life and simulate the spread of a disease (Section 2). We then show how this can generate vast quantities of LBSN data (Section 3) and how this simulation can be used to not only for disease spread but also as a test-bed for various scenarios (Section 4). Section 5 concludes the article with an outlook into challenges that this line of inquiry poses.

2 Geo-Social Agent-Based Simulation and Disease Model

At the core of our approach to generating very large, high fidelity, and socially plausible LBSN data is a novel geo-social agent-based simulation framework [13, 15]. We use the term geo-social to refer to a model that makes explicit use of geographical information and social behaviors in the simulation process. The model simulates individuals (i.e., agents) who live, travel, and interact in an urban setting.

Agents exhibit plausible social behavior that is based on well-respected psychological and social science theories (e.g., [2, 20]). The agents have comprehensive needs and behaviors which result in complex PoL, and temporal social networks. In addition, each simulation instance of our framework is based on (real or synthetic) spatial networks with locations and social networks of users. The agent model logic used to generate the data is constructed based on people’s daily PoL supported by well-respected psychology and social science theories (e.g., [2, 20]). Each individual is equipped with the first three levels of the Maslow’s [20] Hierarchy of Needs: (1) physiological, (2) safety, and (3) belonging and love needs.



Figure 1: Screenshot of the epidemic simulator depicting the French Quarter, New Orleans, LA, USA

When physiological and safety needs are met, an agent may then choose to visit a recreational site (i.e., figurative “hubs”) for the purpose of socialization. At these places, agents may meet new peoples, create budding friendships, and or improve their existing friendship bonds, which are captured using a weighted and directed “social network.” A recreational site visit can sometimes yield a friendship with a stranger if there are no known friends inside (i.e., focal closure). This chance slightly increases if the stranger is actually a friend’s friend (i.e., cyclic closure). Similarly, the lack of social interaction decreases one’s friendship strength, which may eventually lead to break ties. Restaurants are another type of site used by agents to satisfy their physiological needs caused by hunger. Agents who visit restaurants have the chance to meet with their friends without coordination. Such meetings increase the strength of an existing friendship.

The simulation writes spatial and social information for each agent into large, shared log-files which can be processed and analyzed offline. Simulation parameters can be adjusted to create social settings similar to the real-world allowing us to create massive sets of simulated LBSN data. Such high-fidelity data sets contain all individuals of our simulated world with certainty while not impacting the privacy of any human subject in the real-world. As a deliverable, this research yielded synthetic LBSN data sets of thousands of users, scaling to years of observed user data, and thus creating gigabytes of meaningful check-in and social interaction data.

Using the geo-social agent-based model as a basis [13, 15], we have designed a SEIR compartmental disease representation. The SEIR acronym refers a “Susceptible” individual who has the potential to be “Exposed” to a disease. Once exposed, people are potentially “Infected”. Finally, infected individuals transition to a “Recovered” state. For each unique disease, the entire population is initially assumed susceptible. To test different characteristics of diseases, we not only vary exposure and infection times but also person-to-person transmission rates. By doing this, it becomes possible to study many diseases ranging from those that quickly disappear to those that infect a large part of the population. We introduce restaurants as a place where new diseases are initiated and susceptible individuals are exposed based on an environmental spread probability. This spread probability decays daily assuming that the disease pathogen weakens and disappears over time. After being exposed, individuals do not immediately spread the disease but they have to become infectious first.

Figure 1 provides a snapshot of a disease simulation using the geo-social agent-based model [13]. The simulation covers the French Quarter, New Orleans, LA (spatial network) and agent locations are shown on the left in the figure. The current social network is depicted on the right, where nodes correspond to agents and links correspond to friendship. In both networks, the color of an agent corresponds to their disease status. The time series of the number of infections in the center of Figure 1 shows that we have just passed the epidemic peak—similar to the COVID-19 situation now. A full one-minute video simulating two months of simulation time in can be found at our project page [11]. This simulation also generates a large data set (several GBs) that captures simulation parameters and results with a 5min temporal resolution.

Table 1: Data Sets Resulting from Location-Based Social Network Simulation

Settings	# of Users	# of Locations	# of CheckIns	# of Links	Period (month)
GMU-1K	874	113M	2,082,788	9,114,337	15
GMU-1K	874	913M	16,210,909	75,747,439	121
GMU-3K	2,589	335M	6,229,293	27,650,685	15
GMU-5K	4,648	602M	11,189,377	54,250,961	15
NOLA-1K	863	111M	2,099,867	9,160,459	15
NOLA-1K	863	1,647M	29,597,885	141,425,945	221
NOLA-3K	2,720	352M	6,886,573	27,284,999	15
NOLA-5K	4,728	612M	12,007,415	48,710,881	15
TownS-1K	876	113M	2,101,620	7,643,374	15
TownS-3K	2,645	342M	6,454,785	26,364,057	15
TownS-5K	4,349	563M	10,760,008	45,118,825	15
TownL-1K	853	110M	2,030,688	6,418,473	15
TownL-3K	2,550	330M	6,340,360	22,655,915	15
TownL-5K	4,216	546M	10,548,956	40,431,579	15

3 Simulations and Data

To demonstrate the feasibility of generating LBSN data over time, Table 1 gives an overview of the generated output data from the location-based social network simulation [13]. All data sets are available at the Open Science Framework (OSF) [12]. Each data set can be downloaded directly. For low bandwidth connections, a pre-compiled executable of each simulation can be downloaded to re-generated the data locally¹. The table shows the number of agent check-ins and the number of social links attributed to each of scenarios that consist of different numbers (1K, 3K, and 5K) of agents and four different maps. We note that the number of social links may be larger than the square of the number of users. That’s due to the temporal nature of the data set: Social links may emerge and break over time as described in Section 2: Agents meet new friends, but slowly forget about them if their friendship is not reinforced with further meetings. Thus, each link comes with a start time-stamp and end time-stamp.

The four maps are (i) New Orleans, LA (NOLA), (ii) the George Mason University (GMU) campus, (iii) a small synthetic town (TownS) and (iv) a large synthetic town (TownL). The synthetic maps were created using a spatial network and place generator based in a generative grammar similar to L-systems described in [14]. A more detailed description of these data sets can be found in [13].

Figure 2 shows four visualizations of the social networks of 1K agents exemplary for GMU, NOLA, TownS, and TownL at the end of the 15 months simulation. These visuals show different types of network structures, such as two to three large social communities for the synthetic TownL, and one large community for GMU and NOLA. Since it is hard to describe the evolution of a social network over time, we have created a video for each of the four spatial areas showing the social network evolution over the 129,600 steps within the 15 months simulation time. These videos can be found at the project page², and show how the social networks evolve from small isolated cliques into a large and complex network showing different sub-structures.

We welcome the SIGSPATIAL community to utilize the large and high fidelity temporal location-based social network data sets for their research. The purpose of these data sets is not to replace real-world data sets such as Gowalla [6], but to complete real-world data sets with large-scale high-density data to answer the question: How well would algorithms work on large, dense, and ground-truth enriched LBSN data.

¹Source code and executable: <https://github.com/gmuggs/pol>

²LBSN Data Generation project page: <https://mdm2020.joonseok.org>

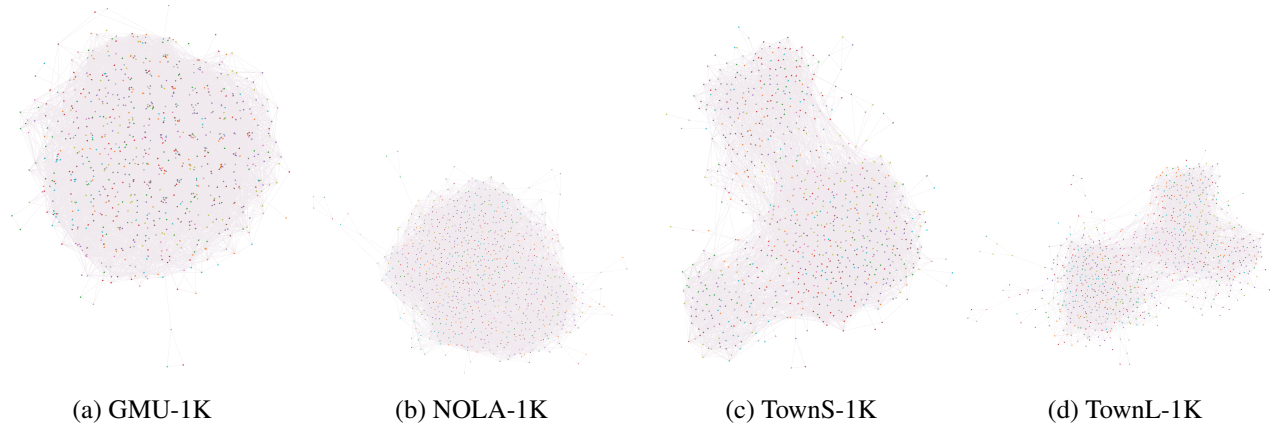


Figure 2: Social network (Note: the location of a node does not represent the location in the spatial network)

Using the Epidemic Model of Section 2 to extend our geo-social agent-based model we can experiment with different types of diseases and disease dynamics. To assess, for example, how epidemiological parameters (e.g., reproduction number) influence disease dynamics, we ran a simulation of two years with randomized epidemiological parameters. Figure 3 shows the number of infectious agents and recovered agents with each (synthetic) disease over time. All diseases are identified as d-ID (e.g., d-12), and the smaller ID number outbreaks earlier than the larger ID number. Only 36 out of 52 diseases outbreak due to epidemiological parameters and diverse environmental aspects. For instance, a disease source might have a very low chance to transmit a disease from environments. It is plausible that some areas that are a potential source of a disease (e.g., d-11) are not well visited. Some diseases (e.g., d-10, d-39) are very contagious and spread quickly during a short period. They have high epidemic peaks. We notice that the numbers of recovered agents with most diseases nearly reach the population, i.e., 5,000 agents. This is because in this scenario no policies are prescribed to mitigate epidemics. Thus, the population continues to mingle as the disease spreads unmitigated. Data for this simulation, including temporal social network data, agent check-in data, and disease data (including exposure, infection, and recovery times for each agent) is available at OSF [16].

4 Geo-Social Simulation and Disease Spread Mitigation

A geo-social simulation is a powerful tool to explore “what if?” scenarios such as in the case of assessing different mitigation measures to limit disease spread. Using agent-based modeling as part of a predictive analytics approach allows us to explore how a change to the simulation (often called an intervention or prescription) will affect the system. Beyond purely predictive analytics, geosimulation [4] allows to explore the parameter space of possible prescriptions to find optimal strategies (or policies) to achieve a desired system state and outcome. We can refer to such a search for optimal policies as *prescriptive analytics*. For readers wishing to learn more about prescriptive analytics please see the 1st ACM KDD Workshop on Prescriptive Analytics for the Physical World (PAPW 2020) [1].

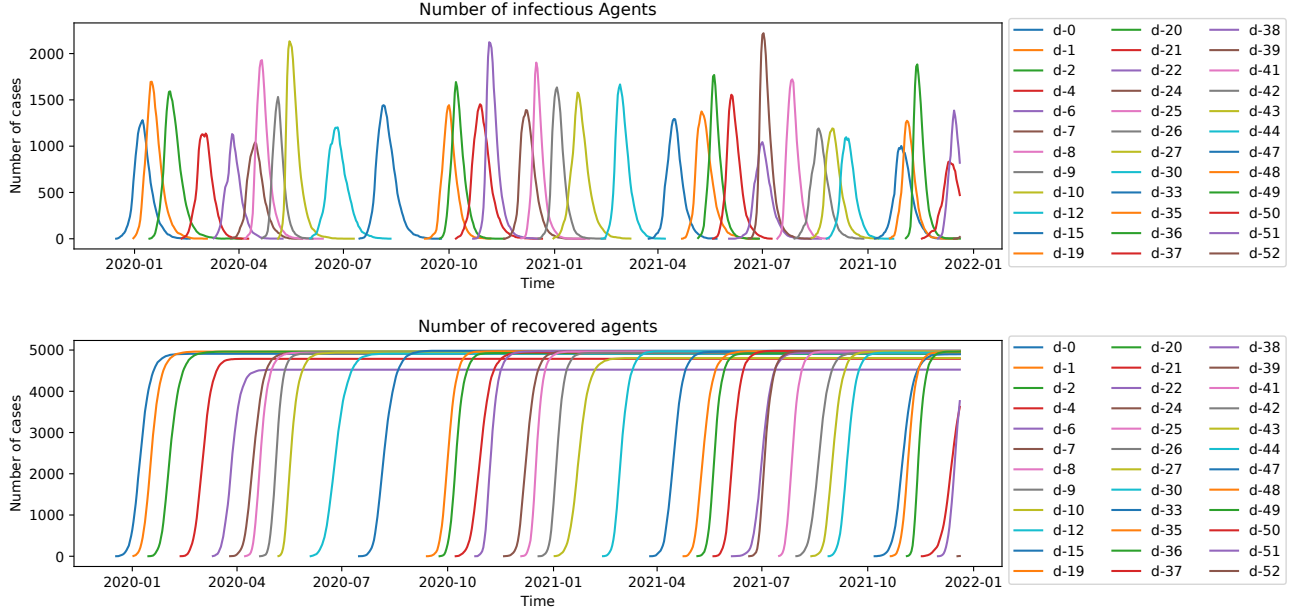


Figure 3: Disease dynamics varying epidemiological parameters

Initial Prescriptive Experiments

We showcase how to leverage our geo-social simulation for prescriptive analytics using two prescribed policies to mitigate the result of a disease. The first policy requires all agents to wear simulated Personal Protective Equipment (PPE) that reduce the chance of infection by 50%. The second policy enforces strict social distancing measures onto a fixed proportion of 50% of the population. Those who follow the social distancing order avoid recreational site visits from meeting people although they still go to restaurants. As a baseline, we also ran a “null-prescription” in which no intervention was prescribed. Figure 4a shows the number of new disease cases immediately after the prescription at a simulation date of 2020-01-01 in all three cases. To quantify the uncertainty of the simulation results, we repeated each of the three scenarios 30 times (differing only in random seeds) and charted the resulting confidence intervals (standard deviation) of new cases. First, we observe that the social distancing prescription was extremely effective. The peak of the curve has been flattened from close to 350 new daily cases to less than 200. However, our simulation shows that merely wearing protective gear without any change in behavior has no significant effect (for the case of this disease). The number of new infections when using protective gear is nearly the same in the case of the null prescription (i.e., take no action). Figure 4a shows a weekly periodicity. This is due to most of our agents not working during the weekend, thus having more time to mingle with others at recreational sites.

Figures 4b, 4c, and 4d show the time series of different disease states of agents for different policies ranging from the null prescription (No Action), the use of protective gear (PPE), and social distancing (Social Distancing Order), respectively. Initially almost all of the 5000 agents are susceptible. As time passes, agents become exposed, with exposure peaks on weekends. Exposed agents become infectious after a while, and thus are able to expose other agents. In the case of No Action and PPE, we see that the number of recovered agents approaches 5000, implying that almost all agents had the disease at one point during these two months of simulation time. In contrast, Figure 4d shows that more than 1000 agents remain susceptible by enforcing social distancing, thus have never been infected even after the disease has disappeared (once all agents are either susceptible or recovered).

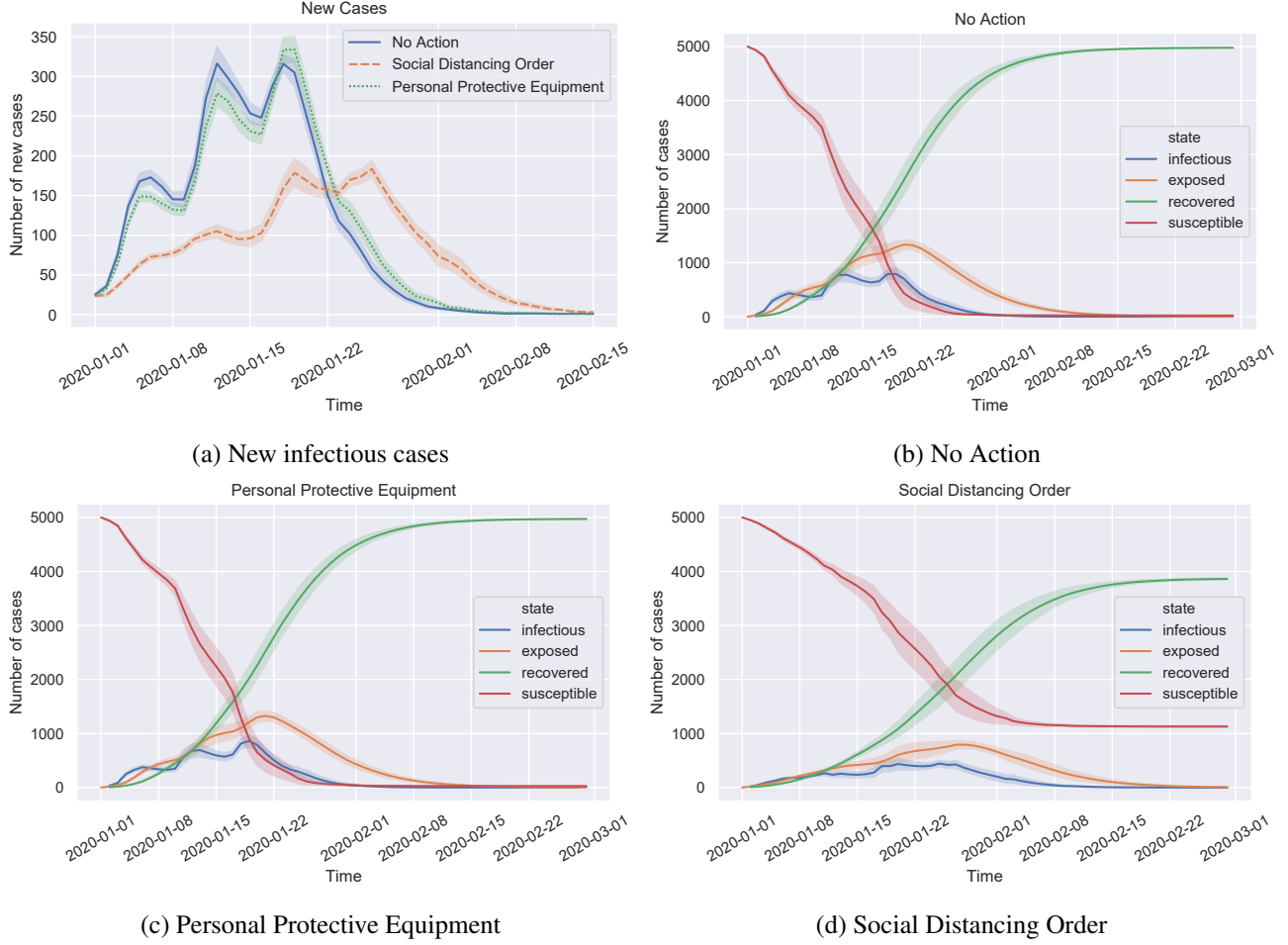


Figure 4: New cases and SEIR epidemic course

Building on these policies, an open question is how to find a policy that minimizes new infections but also minimizes the socio-economic cost of interventions. Not only how many, but which agents should be forced to carry out social distancing? How many, and which agents should be quarantined, and when? Should recreational sites, restaurants and other sites be closed? Which ones and for how long? Answering these questions and finding optimal prescriptions among the combinatorial parameter space is the major challenge of prescriptive analytics.

It is worth noting that we captured four diseases moving through our simulated world in these experiments. Each disease was in a unique phase of its temporal and inter-personal transition across our synthetic, social population and a dynamically emerging social network. Those phases were: 1) a totally new and initial presentation instance, 2) a young (two-week old) instance, 3) a mature disease in full spread, and 4) a near-to-wane disease about to disappear in the population. Further information about experiments can be found at the project website [11]. We expect such a feature will be capable of simulating a simultaneous flu and second-wave COVID-19 epidemic outbreaks recently mentioned as a potential scenario by the CDC director [23].

Our simulation uniquely captures locations and temporal social networks using realistic rules of social behavior based on PoL. As such, we hope that the data sets that we generated, including check-in data, temporal social networks and disease information will help researchers to investigate to what degree it is possible to predict the effect of diseases and social distancing measures on social networks and social needs. With our data providing ground truth data on disease cases, the data also helps investigate the possibility of detecting disease cases and their stages by leveraging both temporal social network and check-in data [16].

5 Conclusions and Challenges

In this SIGSPATIAL Special Newsletter article, we spotlighted our geo-social simulation framework in the context of two applications:

1. To generate very large location-based temporal social network data that captures the temporal change of social networks based on agent behavior grounded in basic social theory. We hope that our data sets that we made available on OSF [12] will be useful in future LBSN research.
2. We leveraged our simulation to simulate the spread of diseases across space, time, and social networks. The generated data sets include check-in data, social networks and disease information. They are available at OSF [16] and should help researchers link LBSN and disease data to answer questions such as: How will social distancing measures disrupt our social fabric? To what degree can we detect and predict new disease cases by leveraging both temporal social network and check-in data?

Looking ahead, one of our bigger challenges is to leverage geo-social simulations for prescriptive analytics in order to find optimal policies to solve a broad range of problems such as traffic congestion and disease spread. For example, what is the best way to leverage work from home policies, and how can these policies be best scheduled among the population in order to tackle the traffic congestion across the world? Or with respect to disease models, what are the best social distancing policies that most effectively contain a disease while also minimizing social and economic damage? But this challenge is not trivial: the space of possible policies is combinatorially large, including parameters for each individual agent (whom to intervene), and for each simulation day (when to intervene). To identify good heuristics to explore this space, we want to apply simulation calibration tools, such as for example dynamic data assimilation [8, 10, 26].

References

- [1] 1st ACM KDD Workshop on Prescriptive Analytics for the Physical World. <https://prescriptive-analytics.github.io/> (accessed 2020-05-16).
- [2] I. Ajzen. The Theory of Planned Behavior. *Organizational Behavior and Human Decision Processes*, 50(2):179–211, 1991.
- [3] S. Bansal, G. Chowell, L. Simonsen, A. Vespignani, and C. Viboud. Big Data for Infectious Disease Surveillance and Modeling. *The Journal of Infectious Diseases*, 214(suppl.4):S375–S379, 2016.
- [4] I. Benenson and P. M. Torrens. *Geosimulation: Automata-Based Modelling of Urban Phenomena*. John Wiley & Sons, London, UK, 2004.
- [5] M. Bithell, J. Brasington, and K. Richards. Discrete-element, Individual-based and Agent-based Models: Tools for Interdisciplinary Enquiry in Geography? *Geoforum*, 39(2):625–642, 2008.
- [6] E. Cho, S. A. Myers, and J. Leskovec. Friendship and Mobility: User Movement in Location-based Social Networks. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1082–1090, 2011.
- [7] A. Crooks, N. Malleson, E. Manley, and A. Heppenstall. *Agent-Based Modelling and Geographical Information Systems: A Practical Primer*. Sage, London, UK, 2019.
- [8] M. D’Auria, E. O. Scott, R. S. Lather, J. Hilty, and S. Luke. Assisted Parameter and Behavior Calibration in Agent-based Models with Distributed Optimization. In *International Conference on Practical Applications of Agents and Multi-Agent Systems (PAAMS’20) (to appear)*, 2020.
- [9] H. Kavak, J.-S. Kim, A. Crooks, D. Pfoser, C. Wenk, and A. Züfle. Location-based Social Simulation. In *Proceedings of the 16th International Symposium on Spatial and Temporal Databases*, pages 218–221, 2019.

- [10] L.-M. Kieu, N. Malleson, and A. Heppenstall. Dealing with Uncertainty in Agent-based Models for Short-term Predictions. *Royal Society Open Science*, 7(1):191074, 2020.
- [11] J.-S. Kim, H. Jin, H. Kavak, O. C. Rouly, A. Crooks, D. Pfoser, C. Wenk, and A. Züfle. Geo-social Simulation Project Website. <https://geosocial.joonseok.org> (accessed 2020-05-16).
- [12] J.-S. Kim, H. Jin, H. Kavak, O. C. Rouly, A. Crooks, D. Pfoser, C. Wenk, and A. Züfle. LBSN-data. <https://osf.io/e24th> (accessed 2020-05-16).
- [13] J.-S. Kim, H. Jin, H. Kavak, O. C. Rouly, A. Crooks, D. Pfoser, C. Wenk, and A. Züfle. Location-based Social Network Data Generation Based on Patterns of Life. In *IEEE International Conference on Mobile Data Management (MDM'20) (to appear)*. IEEE, 2020.
- [14] J.-S. Kim, H. Kavak, and A. Crooks. Procedural City Generation Beyond Game Development. *SIGSPATIAL Special*, 10(2):34–41, 2018.
- [15] J.-S. Kim, H. Kavak, U. Manzoor, and A. Züfle. Advancing Simulation Experimentation Capabilities with Runtime Interventions. In *2019 Spring Simulation Conference (SpringSim)*, pages 1–11. IEEE, 2019.
- [16] J.-S. Kim, H. Kavak, O. C. Rouly, H. Jin, A. Crooks, D. Pfoser, C. Wenk, and A. Züfle. LBSN-disease-data. <https://osf.io/k8qjb> (accessed 2020-05-16).
- [17] M. Li, R. Westerholt, H. Fan, and A. Zipf. Assessing Spatiotemporal Predictability of LBSN: A Case Study of Three Foursquare Datasets. *GeoInformatica*, 22(3):541–561, 2018.
- [18] Y. Liu, T.-A. N. Pham, G. Cong, and Q. Yuan. An Experimental Evaluation of Point-of-interest Recommendation in Location-based Social Networks. *Proceedings of the VLDB Endowment*, 10(10):1010–1021, 2017.
- [19] P. A. Longley, M. F. Goodchild, D. J. Maguire, and D. W. Rhind. *Geographical Information Systems and Science*. John Wiley & Sons, New York, NY, 3rd edition edition, 2010.
- [20] A. H. Maslow. A Theory of Human Motivation. *Psychological Review*, 50(4):370, 1943.
- [21] S. Scellato, A. Noulas, and C. Mascolo. Exploiting Place Features in Link Prediction on Location-based Social Networks. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1046–1054, 2011.
- [22] J. Snow. *On Mode of Communication of Cholera*. Churchill, London, England, 1855.
- [23] L. H. Sun. CDC Director Warns Second Wave of Coronavirus is Likely to be Even More Devastating. <https://www.washingtonpost.com/health/2020/04/21/coronavirus-secondwave-cdcdirector/> (accessed 2020-04-21).
- [24] M. M. Waldrop. News Feature: Special Agents offer Modeling Upgrade. *Proceedings of the National Academy of Sciences*, 114(28):7176–7179, 2017.
- [25] H. Wang, M. Terrovitis, and N. Mamoulis. Location recommendation in location-based social networks using user check-in data. In *Proceedings of the 21st ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, pages 374–383, 2013.
- [26] J. A. Ward, A. J. Evans, and N. Malleson. Dynamic Calibration of Agent-based Models using Data Assimilation. *Open Science*, 3(4):150703, 2016.
- [27] M. Ye, P. Yin, and W.-C. Lee. Location Recommendation for Location-based Social Networks. In *Proceedings of the 18th SIGSPATIAL International Conference on Advances in Geographic Information Systems*, pages 458–461, 2010.
- [28] J.-D. Zhang and C.-Y. Chow. Point-of-interest Recommendations in Location-based Social Networks. *SIGSPATIAL Special*, 7(3):26–33, 2016.