# Using Interacting Particle Systems to Measure Spread of Misinformation

Andrew Wang, Eero Saar Gallano, Kenny Kang

## I. INTRODUCTION

The permeation of "fake news" and misinformation is one of the greatest obstacles towards maintaining a functioning democracy. We model and measure the spread of misinformation against adjustments to parameters of the model to explore outcomes and preventative measures.

Interactive particle systems have been used to model diseases and political leanings in the context of networks, social or physical. We extend the utility of interactive particle systems to modeling spread of misinformation within social networks by utilizing various aspects of different IPSs. While the nature of IPSs are stochastic, their dependency on a graph structure permits the replication of complexity of large-scale systems. Furthermore, we can run numerous simulations of varying scenarios to find approximations of the underlying patterns inherit to the system we are attempting to model. For these reasons, IPSs are an ideal candidate for modelling the spread of misinformation in social networks.

## II. HYPOTHESIS

With respect to the $n$ being the number of non-source nodes in a graph, we posit that within $\log(n)$ time steps, at least $\frac{1}{2}n$ nodes will be affected by misinformation. Simulations with multiple origins will reach majority infection on average faster than a single source, but does not necessarily scale linearly.

## III. MODEL

The interactive particle system will employ two submodules that will work in tandem: a traffic model to regulate the flow of information through a node and a general voter model to simulate the interactions of social network groups. More specifically, the misinformation model will contain three object types: sources, nodes, and packets.

(a) Packets are mutable instances of information that propagate throughout the network.
(b) Nodes serve as intermediate points of information flow through the network where packets are copied and mutated.
(c) Sources are points in a network that generate new, unmutated packets.

### A. Packet

Packets will have two properties, mutation and time stamp. Mutation is the amount a packet deviates from its original form whereas the time stamp is the time the original packet was generated. Packets are distinguished only by their original form. Hence two packets with different levels of mutation that originate from same unmutated packet will not be differentiated, and thus a combined form representing these two packets will be sent. For our particular model, we simply sum the mutation levels of duplicate packets, though other realistic operations such as averaging can be considered.

### B. Node

Nodes will have two properties, propagation parameter $\lambda_v$, mutation probability $\rho_v$, mutation distribution $Q_v$. The propagation parameter generates a distribution $\text{Exp}(1/\lambda_v)$ that determines the time at which the next packet is sent. The mutation probability is a parameter to determine whether a mutation occurs or not ($\text{Bernoulli}(\rho_v)$). The mutation distribution is a continuous distribution of mutation values. Note that normal distributions are not expected to represent the differences between information learned and the (mis)information shared between nodes though they very well can.

### C. Source

Sources, akin to nodes, will have a send rate $\lambda_v$ as well. However, a key distinction is that sources are unable to receive or mutate packets.

### D. Network

Though not an object in a model, how the network is assembled or generated is an important topic of discussion. While Erdős–Rényi model is the canonical random graph, we instead opted for scale-free networks as they seem to more closely resemble social networks, though this remains a point of contention.[1]

## IV. EXPERIMENTS

As our hypothesis is concerned with the affects of number of sources versus the vulnerability of individuals to misinformation, the only varying hyperparameter considered is the number of source points in the network. All other parameters remain fixed. Each experiment runs for 250 ticks, with 250 non-source nodes in each graph. Each node parameter is i.i.d. sampled from the following distributions: $\lambda_v \sim |\mathcal{N}(0, \frac{1}{100})|$, $\rho_v \sim \text{Uniform}[0, 1]$, and $Q_v = \mathcal{N}(0.5, 1)$. Furthermore, the random graphs are generated from a free-scale generator algorithm with parameters, $\alpha = 0.48$, $\beta = 0.04$, $\gamma = 0.48$. We range the number of sources from $n_{\text{sources}} \in \{1, \ldots, 25\}$.

---

[1]Discussion is based on the paper by B. Bollobás, C. Borgs, J.Chayes, and O. Riordan, Directed scale-free graphs, Proceedings of the fourteenth annual ACM-SIAM Symposium on Discrete Algorithms, 132–139, 2003.

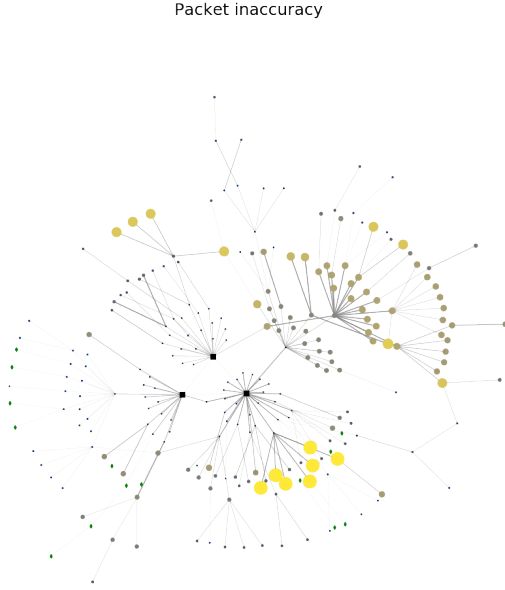A sample simulation with the parameters noted above produce the following metrics [2]:

Packet inaccuracy



Fig. 1. Graph visualization of cumulative packet inaccuracy
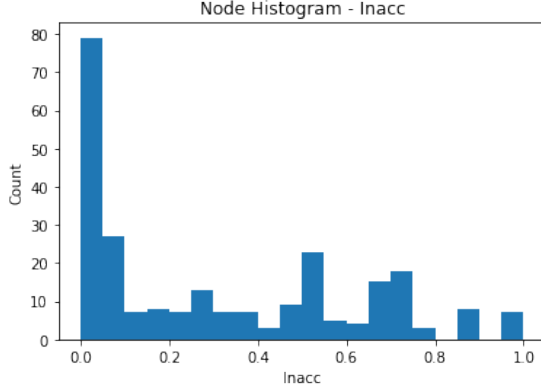


Fig. 2. Histograph visualizing distribution of cumulative packet inaccuracy

Graph visualizations help illustrate a holistic picture of the simulation run; node size and color indicate the magnitude of the metric present in that node. Other, subtle statistics are also revealed in the graph by the thickness of each edge, indicative of high or low packet traffic, and the shape of the node, where green diamonds represent nodes that were not exposed to misinformation. We can furthermore visualize the distribution of these metrics via a histogram. Other observable metrics include packet recency (how quickly a packet reaches a node) and packet flow-through (how many packets are sent *through* a node).

---

[2]Note that the code accompanying this paper is not fully described. Utility methods and attributes are included to allow cumulative measures across the span of the simulation. Code here: https://github.com/gallanoe/mis-sim

## V. CHALLENGES

In order to have a model that reflect the structure and distribution of connections between people, our graph network needs to have a support a specific form of construction and a large population size. Simulating and visualizing large graph networks is computationally expensive and messy, but with the aid of various open source libraries, we were able to construct our simulation and run it within reasonable time.

The $50\%$ by $\log(n)$ hypothesis was not being reached by the exact number of time steps as we had guessed. We added an alternative hypothesis of $2\log(n)$ as our collected data had shown properties indicating phenomena similar to our hypothesis. Confirmation of this alternate hypothesis would imply our original hypothesis was correct by asymptotic equality assumptions.

## VI. OBSERVATIONS

Averaging over twenty simulation runs for each value of $n_{\text{sources}}$, we observe drastic differences between the amount of time between exposure to information versus exposure to misinformation. Note the $50\%$ y-axis markers and $\log(n)$ and $2\log(n)$ x-axis markers.
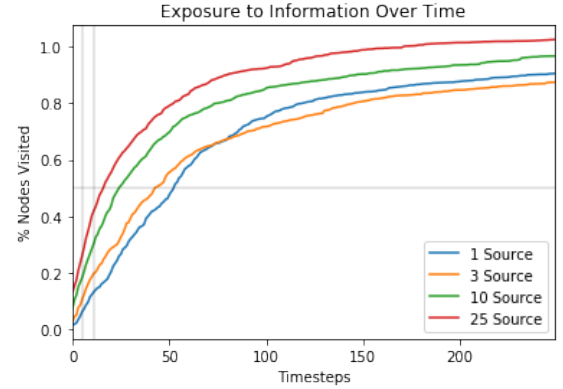


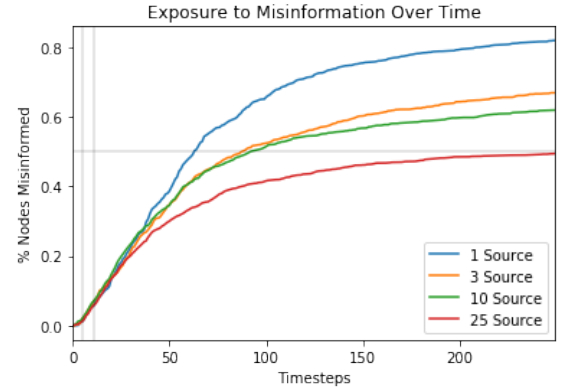Fig. 3. % nodes who are exposed to information over time



Fig. 4. % nodes who are exposed to misinformation over time

Note how appropriate our hypothesis seems to be when applied to information-exposed nodes but less so when measuring misinformation exposure. Furthermore, while there is

an obvious asymptotic bound at $y = 1$, this asymptotic bound for misinformation decreases as the number of sources increase.

However, it must be noted that nodes that are not exposed to information at all are not the topic of concern. As they have had no exposure to information, by definition they cannot be exposed to misinformation, while the ones that have been exposed to information are also possibly exposed to misinformation. Though running the simulation further may eventually lead to exposure to information for these non-exposed nodes, it is arguable to say that they are irrelevant to the hypothesis. As such, we also include examination of ratios between misinformation-exposed nodes and information-exposed nodes.
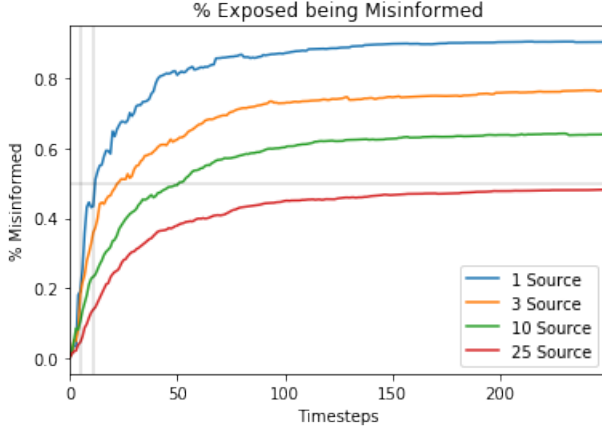


Fig. 5. % of information-exposed nodes who are misinformed over time

With this adjustment to metrics, our hypothesis becomes partially supported, namely for $n_{sources} = 1$. However, similar observations compared to those gathered from previous metric comparisons can be made here.

Finally, to affirm the relationship between $n_{sources}$ and % misinformed, we plot it against the average number of time steps needed to expose $50\%$ of nodes to misinformation.



Fig. 6. Avg. time to achieve $50\%$ misinformation-exposure

Misinformation slows down drastically as the number of sources increase, while the dissemination of information, misinformed or not, increases. Note that when $n_{sources} = 25$, the minimal time step for majority misinformation exposure dips to $t = 0$. This is simply due to the fact that with that number of sources, we simply did not achieve the majority within the 250 time steps allocated for each simulation run. This strengthens the hypothesis that misinformation is slowed down when the number of sources increase. In fact, even with such limited data, we can suggest that the number of time steps needed to reach misinformed majority is at least linear with respect to the number of sources (though evidence arguable suggests is grows more so).

## VII. Conclusion

In conclusion, we found that for our experimental parameters, we did not reach $50\%$ misinformation within $\log(n)$ time steps for any of our source node counts. However, when only considering nodes that have been exposed to information, our revised hypothesis is affirmed for only one source of information. Any more and the dissemination of misinformation becomes negated. In addition, the number of time steps it takes to achieve a majority misinformed population grows at least linearly with respect to the number of sources.

## VIII. Future Work

However, these experiments are executed under specific assumptios. When contextualized with our misinformation modeling problem, we assume that our sources are perfectly unbiased as the source nodes do not generate packets with any initial mutation. Furthermore, we did not make any differentiation between different types of nodes; for example, a newspaper might have a higher rate of packet mutation than that of a regular person who's having friendly discussions with their peers.

Though much discussion is still to be had about the modeling of social networks with scale-free graphs, the inherit structure of scale-free graphs allows us to classify different types of nodes from another. It's possible to assign the nodes their appropriate roles and optimize the parameters for generating scale-free graphs to more closely resemble social networks of the real world. In addition, we can explore the distribution of parameters/attributes for each nodes.

Finally, the flexibility in notation/terminology allows us to employ this model to other fields such as disease exposure and mutation.