

1 **Reviewer #2**

2 – Generic NN formulation and relation to RNN: Our (14) is a highly structured neural network (NN) due to the
3 Mori-Zwanzig moment closure (please note the form of \mathbf{f} in Eq.(7)). This enables us to learn the diffusion parameter
4 and estimate influence (both for the diffusion network) using the mean-field dynamical system (14). This is different
5 from learning a generic RNN which is not interpretable and does not serve the purpose of learning diffusion parameters.

6 – Relevance to NeurIPS: Learning diffusion parameters and predicting influence are one of central problems in the ML
7 community. Please see [14-17],[21-23].

8 – Input is only the initial condition: This is because we need to estimate the influence of source set during testing phase,
9 where no additional temporal data is available. Hence we can only use the source set as input of our network during
10 training and testing. This is different from the standard RNN applications.

11 – Map to standard RNN/LSTM: Yes, we can map a dynamical system to other types of RNNs. However, our purpose is
12 to derive a structured dynamical system for influence prediction, and need to integrate the diffusion parameters in the
13 dynamical system for training in this work.

14 – A more precise definition of ε : in the present work, we set it to a generic NN to approximate the MZ memory kernel.
15 We will state this with more details if a revision is allowed.

16 **Reviewer #4**

17 – Splitting the linear and nonlinear part in ε : This is because the linear part is a consequence of the kernel approximation
18 of \mathbf{h} in Eq.(10) and \mathbf{K} in Theorem 3. In this case, we can interpret the matrices \mathbf{B} and \mathbf{C} (and impose proper
19 regularization or prior information if possible). The remainder part of ε is set to generic NN, but not interpretable
20 anymore.

21 – Connection to PMP: Unlike the standard optimal control and Ref.[37], our formulation does not allow time-varying
22 control variables. This requires a modification of the PMP to interpret our network training. By introducing the total
23 Hamiltonian, we showed that the standard back-propagation is directly equivalent to maximizing the total Hamiltonian.
24 This is different from Ref.[37] which relies on successive approximation (different from back-propagation) and is
25 computationally more expensive,

26 – Error of NMF decrease as time increase: this is because as time goes, all nodes will eventually be infected (if the
27 diffusion network is connected). Therefore, all methods (not only NMF) will have lower prediction error as time goes
28 to infinity. But NMF appears to be more accurate in predicting the error in early to middle stage of the propagation
29 (rather than asymptotical error), which is of most interests in practice (e.g., one would like to predict the spread of news
30 in three days rather than a month).

31 – $e(t)$ is not defined: $e(t)$ is defined in Eq.(5) above Theorem 1.

32 **Reviewer #5**

33 – Evaluation of approximation error: please note that it is *infeasible* in practice to evaluate the original MZ memory term,
34 because exact evaluation involves solving an ODE system of size 2^n (n is the size of the diffusion network). Instead,
35 we use MC simulations to approximate ground truth as reference, and compared our results against it, as shown in our
36 experiment part. This has been the standard approach for comparison in the literature.

37 – Using InfluoLearner in influence maximization: Our comparison with InfluoLearner on influence prediction demonstrated
38 the significant improvement in accuracy. Therefore its performance cannot exceed NMF in influence maximization
39 which heavily relies on the accuracy of influence prediction. Moreover, the computation complexity of InfluoLearner is
40 much higher than NMF, and not suitable for large scale problems and dense networks in influence maximization.

41 – Pseudo-code: the NMF network structure is given in Eq.(14), which can be solved by standard network training. We
42 will provide more details if a revision is allowed.

43 – MAE only on 100 source sets: the 100 source sets refers to the testing data, whereas 1,000 refers to the training data.

44 **Reviewer #6**

45 – Examples: the network platform (like Twitter), campaign company, or advertiser are often interested in the influence
46 of individuals (users) on networks—if the individual posts a news/advertisement, then it will be seen and retweeted by
47 his/her followers, and then the followers of them, and so on. Accurate influence estimation proposed in our work helps
48 the campaign company or advertiser to decide which individuals to select, at a certain monetary cost, to spread the news
49 as fast as possible (to maximal amount of people on the network). It also has applications in disease intervention where