# Parallel Machine Learning for Forecasting the Dynamics of Complex Networks

Keshav Srinivasan,[1] Nolan Coble,[1, 2] Joy Hamlin,[3] Thomas Antonsen,[1] Edward Ott,[1] and Michelle Girvan[1]

[1] *University of Maryland, College Park, Maryland 20742*
[2] *SUNY Brockport, New York 14420, USA.*
[3] *Stony Brook University, New York 11794, USA*
(Dated: August 30, 2021)

Forecasting the dynamics of large complex networks from previous time-series data is important in a wide range of contexts. Here we present a machine learning scheme for this task using a parallel architecture that mimics the topology of the network of interest. We demonstrate the utility and scalability of our method implemented using reservoir computing on a chaotic network of oscillators. Two levels of prior knowledge are considered: (i) the network links are known; and (ii) the network links are unknown and inferred via a data-driven approach to approximately optimize prediction.

Machine learning (ML) has played a vital role in recent scientific advances in many disciplines. A key problem in these contexts is time series prediction of a dynamical system for which a first-principles, knowledge-based description is unavailable [1]. By using ML in combination with measured time-series data, one can hope to construct a faithful model of a system's dynamics and to then use this model to predict the future evolution of the system's state. Our aim in this paper is to address this goal for *large systems of interacting components with complex connectivity and dynamics* - a system type of enormous technical and scientific interest in many fields, ranging, e.g., from neuroscience to power grids. However, straightforward application of the standard ML prediction schemes becomes problematic when applied to forecasting the dynamics of large networks. To deal with such systems, we propose a parallel forecasting method for networks with complex dynamics. In our approach, we construct an ML architecture that mimics the topology of the network. Each node of the network to be predicted is assigned an individual small ML device and these individual ML devices are linked to each other based on the underlying connectivity of the network (either known *a priori* or inferred from the available time series data). We demonstrate and test this method by applying it to a network of Kuramoto oscillators [2, 3] constructed to exhibit chaotic dynamics. Our method is motivated in part by previous work on parallel ML prediction of large spatiotemporally chaotic systems [4, 5].

We consider two scenarios: (a) the connectivity of the oscillator network is known, and (b) the connectivity of the oscillator network is unknown *a priori*, yet may be approximately inferred from node time series data. Scenario (a) serves two purposes: first, as preparation for the more challenging situation presented by scenario (b), and second, as a method applicable to cases where the connectivity is, in fact, known. The main conclusion of our paper is that our proposed parallel ML scheme enables data-based network dynamics prediction in cases that would otherwise (i.e., without parallelization) be unattainable.

In order to demonstrate and test our approach, we con-sider the well-studied Kuramoto model of $N$ network-coupled oscillators,

$$\dot{\theta}_i = \omega_i + K \sum_{j=1}^{N} A_{ij} \sin(\theta_j - \theta_i), \qquad (1)$$

where $\theta_i$ is the phase angle of oscillator $i$, $\omega_i$ is the natural frequency of oscillator $i$ when uncoupled, $K$ is the coupling strength, and $A_{ij}$ is the adjacency matrix that specifies the structure of the oscillator network ($A_{ij} = 1$, if there exists a network link from node $j$ to node $i$ with $i \neq j$, and $A_{ij} = 0$, otherwise). Here we consider an undirected ($A_{ij} = A_{ji}$), frequency assortative Kuramoto network [6]. By 'frequency assortative' we mean that two nodes are more likely to be linked if their natural oscillation frequencies are numerically close. The resulting frequency assortative system has chaotic dynamics for certain choices of parameters [7], hence serving as a good example of complex network dynamics. Each node is taken to have the same number of connections (this number is called the node's degree). The oscillator natural frequencies, $\omega_i$, are drawn from a uniform random distribution from $-\pi/2$ to $\pi/2$. The frequency assortative network (i.e., the set of matrix elements $A_{ij}$) is constructed by starting with $N_o$ unlinked nodes, each with its assigned frequency ($\omega_i$ for node $i$) and then successively adding links, as follows. After randomly choosing a node $i$ that still requires additional links, we next randomly pick another node $j$ (not already connected to node $i$) which also still requires additional links, and then, with probability $p_{ij}$, we link nodes $i$ and $j$, where

$$p_{ij} \propto \frac{\delta^{\gamma}}{\delta^{\gamma} + |\omega_i - \omega_j|^{\gamma}} \qquad (2)$$

We continue in this way to make links until all nodes have the desired degree. Due to the form of $p_{ij}$ (Eq. (2)) nodes with similar natural frequencies are connected with a higher probability (see Supplementary Fig. 1). We use the global order parameter, $R$, as a metric to measure

the dynamics of the oscillator network, where

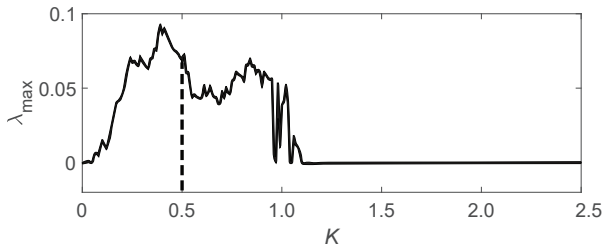$$R(t) = \sum_{i=1}^{N} \sum_{j=1}^{N} A_{ij} e^{i\theta_j} \qquad (3)$$



Figure 1. Largest Lyapunov Exponent as a function of the coupling constant, $K$. The dashed line represents the chosen value of $K$ for our studies.

For our frequency assortative network, with $N = 50$, a nodal degree of 3, $\delta = 0.8$, $\gamma = 5$ and $K = 0.5$ (our standard parameter set for most of out subsequent numerical experiments), we observe chaotic behavior, which is confirmed by the positive value of $\lambda_{max}$, the largest Lyapunov exponent of the system (Fig. 1).

*Background on non-parallel Reservoir Computing prediction.* In this paper we use Reservoir Computing (RC) [8, 9] as our ML scheme, because of its demonstrated utility for time series prediction [1, 10, 11]. We consider a reservoir computer constructed with an artificial high dimensional dynamical system, known as the reservoir, which is coupled to an input through an input layer, specified by a matrix $\mathbf{W_{in}}$ which maps the input vector, $\mathbf{u}$, at discrete time $t$, to the reservoir state variables, which are collectively expressed as the scalar components of the reservoir state vector $\mathbf{r}$. In our RC implementation, the reservoir is a network (not to be confused with the network, e.g., Eq. (1), whose state we desire to predict), and the $k$th component of the vector $\mathbf{r}$ is the scalar state of reservoir node $k$. The RC network is directed, sparse, and random with $N_r$ nodes having average input degree, $\kappa = 3$. The RC adjacency matrix is denoted $B$, with matrix elements $B_{kk} = 0$, and $B_{kl}$ for $k \neq l$ chosen randomly and uniformly from $[-\beta, \beta]$ where $\beta$ is chosen to yield a maximum eigenvalue of $B$ denoted $\rho$ (known as the spectral radius). Each input to the reservoir is sent to $N_r/N_{in}$ reservoir nodes, where $N_{in}$ is the number of inputs to the RC (Note: $N_r$ is chosen to be an integer multiple of $N_{in}$). The input matrix, $\mathbf{W_{in}}$, is then a $N_r \times N_{in}$ dimensional matrix. The elements of $\mathbf{W_{in}}$ are chosen so that every node in the reservoir receives exactly one input from $\mathbf{u}(t)$ while each input in $\mathbf{u}(t)$ is connected to $N_r/N_{in}$ nodes in the reservoir network (see Supplementary Material for further discussion). The non-zero elements are drawn from a uniform random distribution

from $[-\sigma, \sigma]$, where $\sigma$ is the input scaling. The reservoir state, $\mathbf{r}(t)$, is taken to evolve according to

$$\mathbf{r}(t + \Delta t) = \alpha \mathbf{r}(t) + (1 - \alpha) \tanh[\mathbf{B}\mathbf{r}(t) + \mathbf{W_{in}}\mathbf{u}(t)], \quad (4)$$

where the tanh function is applied component-wise to its vector argument. Here $\alpha$ is the leak rate which controls the timescale of the reservoir nodes. The output of the system, $\tilde{\mathbf{u}}$, is defined through the output layer and is given by

$$\tilde{\mathbf{u}}(t) = \mathbf{W_{out}}\mathbf{r}(t). \qquad (5)$$

For the task of time-series prediction, the reservoir computer is used in two different modes: a training mode and a prediction mode. In the training mode, the reservoir computing system, represented by Eqs. (4) and (5), is run for the time interval over which training data $u(t) = u(n\Delta t)$ $(n = -n_t, (1 - n_t), (2 - n_t), ..., 0)$ is available, $\mathbf{r}(n\Delta t)$ is computed, and the output matrix $\mathbf{W_{out}}$ is adjusted ('trained') so that the output of the reservoir computer $\tilde{\mathbf{u}}(t)$ best approximates $\mathbf{u}(t)$. This is done through a ridge regression procedure, wherein we minimize the error summed over the training times $t = n\Delta t$ for $n$ running from $1 - n_t$ to 0,

$$\min_{\mathbf{W_{out}}} \left\{ \sum \left[ \|\mathbf{W_{out}}\mathbf{r}(t) - \mathbf{u}(t)\|^2 \right] + \beta \operatorname{Tr}\left( \mathbf{W_{out}}\mathbf{W_{out}^T} \right) \right\} \tag{6}$$

Here $\beta$ is the Tikhonov regularization parameter that is used to prevent over-fitting. The quantities $(N_r, \rho, \sigma, \alpha$ and $\beta)$, referred to as hyperparameters of the reservoir computing setup, are collectively used to control the performance of system. In this paper we chose the hyperparameters by a subsequent iterative process approximately maximizing the valid prediction time (See Eq. (8)) over the hyperparamters via a coarse grid search (See Supplementary Material Section III). In the prediction mode, the reservoir state now evolves autonomously in "closed-loop" mode; i.e., the output at time $t$, now serves as the input at time $t + \Delta t$,

$$\mathbf{r}(t + \Delta t) = \tanh[\mathbf{B}\mathbf{r}(t) + \mathbf{W_{in}}\mathbf{W_{out}}\mathbf{r}(t)]. \qquad (7)$$

This procedure generates a predicted time series $\hat{\mathbf{u}}(n\Delta t) = \mathbf{W_{out}}\mathbf{r}(t)$ that is assumed to approximate the true future evolution of the state of the system, $\mathbf{u}(t)$ at a time $n\Delta t$ for $n > 0$ (we choose $\Delta t$ small compared to the time scale for variation of $u$ so that $u(n\Delta t)$ essentially specifies the continuous time function $\mathbf{u}(t)$).

*Parallel ML scheme for network prediction.* In order to address the high computational complexity of predicting large networks, we introduce a parallel network RC architecture (see the schematic in Fig. 2). Each node, $i$, in the predicted network is assigned its own reservoir, $R_i$. The inputs to this reservoir are the signal of node $i$ itself, as well as that of the nearest network neighbors

of node $i$. The number of such neighbors is equal to the network degree. The reservoir $R_i$ is then trained on these inputs to predict the signal of node $i$. Because each $R_i$ predicts just one node, its size $N_r$ can be relatively small. In addition, since our parallel scheme uses an interconnected network of independently trained reservoirs, we can efficiently parallelize our training process, making the system scalable to large networks.
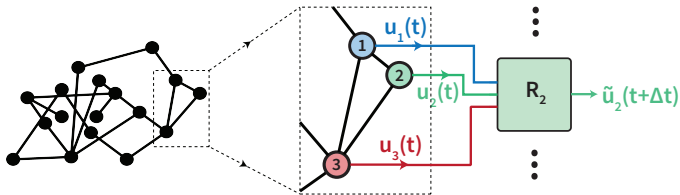


Figure 2. A schematic diagram for the parallel network ML architecture. Here we show Reservoir 2 ($R_2$), which receives input from its assigned node (node 2), plus inputs from nodes connected to node 2 (i.e., nodes 1 and 3). $R_2$ is then is then trained to predict its assigned node (node 2). This process is the same for each node in the network, such that the connectivity among the reservoirs mimics the network to be predicted.

*Results.* To compare the parallel, multiple RC scheme with the single RC approach, we use a $N_o = 50$ node frequency assortative Kuramoto oscillator ($\delta = 0.8$. $\gamma = 5$) network with a coupling constant of $K = 0.5$. We study the magnitude of the global order parameter $|R|$ which tells us about network-level activity (see Fig. 3) and the prediction of the evolution of individual node states (see Supplementary Figure 3), both of which show the same main qualitative behavior. For the purpose of forming inputs to the reservoir, we specify the state of the oscillator $i$ as $[\sin\theta_i(t), \cos\theta_i(t)]$. The input matrix is generated as described above and in Supplementary Material, Section 2.

*Single non-parallel reservoir prediction.* The single reservoir computer prediction can fail as the size of the network we want to forecast increases. This is clearly demonstrated in Fig. 3(a), where the prediction breaks down in a fraction of a Lyapunov time, $\lambda_{max}t$. We quantify the duration of an accurate prediction by a metric that we call the "valid prediction time". This metric is defined as the amount of time elapsed before the normalized root mean squared prediction error (NRMSE), $E(t)$, exceeds some chosen value $f$, $0 < f < 1$, for the the first time, where

$$E(t) = \frac{\|\mathbf{u}(t) - \widetilde{\mathbf{u}}(t)\|}{\langle\|\mathbf{u}(t)\|^2\rangle^{1/2}}. \tag{8}$$

The valid prediction time for $f = 0.1$ is marked in Fig. 3 by a vertical dotted lines. Even for the very large reser-
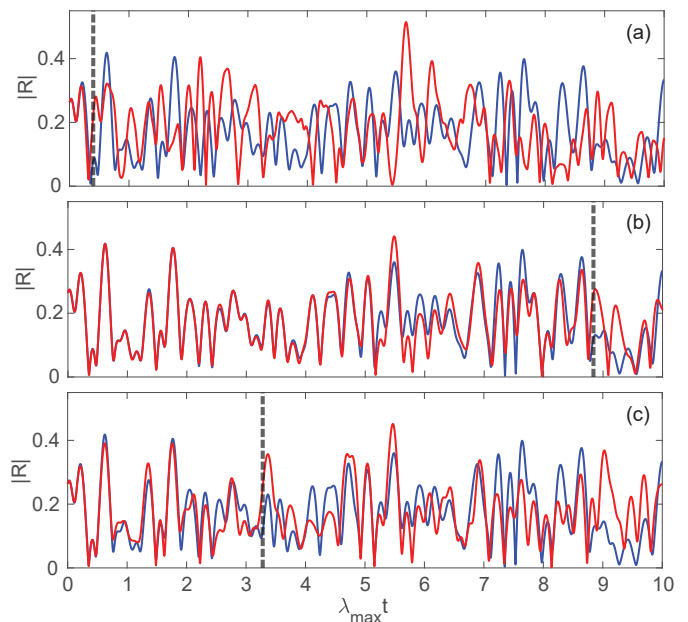


Figure 3. Prediction of the order parameter in 3 different cases of ML implementations. The blue curves are the data (the truth) and red curves are predictions. The dotted lines in each plot denote the valid prediction time. (a) Single, non-parallel RC prediction using a large reservoir ($N_r = 10000$), (b) Parallel prediction with known network links, using 50 separate reservoirs each having modest size ($N_r = 200$), (c) Parallel prediction with unknown network links, using 50 separate reservoirs each having modest size ($N_r = 200$). The network structure is estimated by using transfer entropy as a metric to draw network edges.

voir ($N_r=10000$), close to the limit of our computer resources, that is used in Fig. 3(a), the system is still not able to predict past a fraction of a Lyapunov time.

*Parallel scheme with known links.* In cases where the network structure is known *a priori*, such as in certain social networks, and we can construct our parallel reservoir architecture by using the known network links. In the case of our Kuramoto oscillator network, we demonstrate our results on the 50 node network by using 50 separate parallel reservoirs of relatively modest size ($N_r = 200$ as compared to $N_r = 10000$ for Fig. 3(a)), each having the same set of hyperparameters [see Fig. 3(b)]. The predictive performance of this architecture could potentially be enhanced by individually optimizing hyperparameters for each of the 50 reservoirs, but this would considerably increase both the time and computational resources required for this task. As seen from Fig. 3(b), our parallel scheme does exceedingly well for multiples of the Lyapunov time, $\lambda_{max}t$. This is particularly clear from a comparison of valid prediction time (vertical dashed lines) in Fig. 3(a) versus those in Fig. 3(b), the latter being $\gtrsim 10$ times larger, while at the same time being much less computationally demanding

(mainly due to the difference in $N_r$, $N_r = 10000$ for the nonparallel case versus $N_r = 200$ in the parallel case).

*Parallel scheme with unknown links.* In many cases of interest, one may not have information about the underlying network structure. Using nodal time-series data for finding links in networks, such as metabolic [12] and gene-regulatory networks [13], is an active area of current research. Many heuristic-based [14] and statistics-based tools like conditional mutual information [15], and correlation [16], as well as a machine learning technique [17], have been used for link inference and might give useful approximations of the underlying network structure. These methods could then potentially be used in our parallel network scheme. As an example, we now demonstrate the performance of our parallel method combined with the use of Transfer Entropy [18] for link inference. Transfer entropy is a statistical method to infer causal relationships between variables by using conditional probabilities: If a signal A has a causal effect on signal B, then the probability of B conditioned on its past is different from the probability of B conditioned on both its past and the past of A. Transfer entropy can also be expressed in terms of the conditional mutual information as

$$T_{X \to Y} = I(Y_t; X_{t-1:t-L} | Y_{t-1:t-L}). \tag{9}$$

Considering the problem of network state prediction, we use past measured nodal state time-series to calculate the transfer entropy between each pair of nodes in the network using Eq. (9) and then pick a threshold. Pairs of nodes with transfer entropy values above the threshold are assigned a link between them. As we decrease this threshold, we draw more links and hence increase the average number of supposed neighbors for each node. Initially, decreasing the threshold, or in other words increasing the number of supposed neighbors, increases the number of true positive links and improves the predictive performance of our reservoir scheme. But if this threshold is decreased too much, the number of false positives increases drastically degrading the predictions. Since our goal is prediction, we view the link-inference threshold on the transfer entropy as an additional hyperparameter and chose it (along with the other hyperparameters), so as to optimize the valid prediction time. By this procedure, we effectively bootstrap our prediction process to determine the link threshold criterion. An example set of results for $N_r \approx 200$ (See Supplementary Material Section II for details) is shown in Fig. 3(c). Again, in marked contrast with the results in Fig. 3(a) for a large single RC ($N_r = 10000$), we obtain good predictions, e.g., a valid time between 3 and 4 Lyapunov times for $|R|$.

*Dependence on the size of the predicted network.* Fig. 4 shows a plot of the valid prediction time as a function
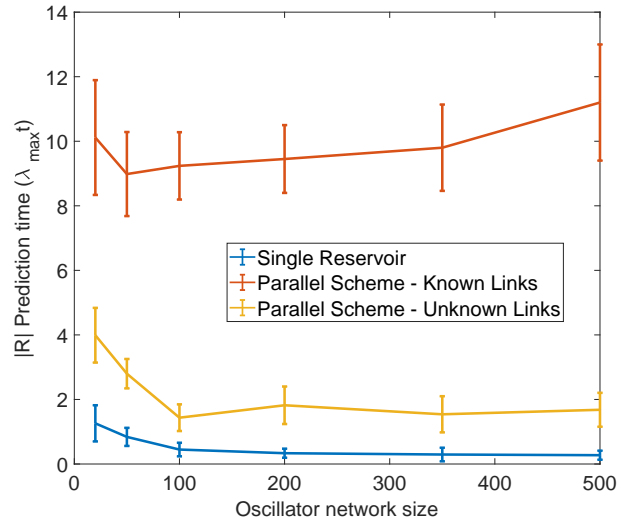


Figure 4. Performance of the different Reservoir Computing methods as a function of the Kuramoto oscillator network size.

of the oscillator network size, $N_o$. As we increase the size of the oscillator network, the prediction using a single reservoir ($N_r = 10000$) quickly degrades even further and becomes unable to capture the network dynamics at all. Since the parallel method assigns a reservoir to each oscillator in the network, for the case with known links, as expected it maintains constant performance to within the estimated uncertainty of the valid times. However, when we do not know the links, the ability of our parallel scheme to predict the network dynamics is limited by its ability to make accurate link predictions. The reduction in performance incurred by missing real links (i.e. false negatives) far outweighs that associated with false positive links. As we increase the oscillator network size, the accuracy of link determination gets worse which ultimately affects the predictive performance of this method. For small oscillator network size of $N_o = 10$, our true positive rate is 100%, while our false discovery rate is 37.5%, but as the network science increases to $N_o = 500$, our true positive rate drops to 96% and our false discovery rate becomes 69.6%. That the prediction degradation is more sensitive to a false negative link inference than to a false positive link inference can be understood as follows. The false negative inference of a link to node $i$ deprives reservoir $R_i$ of vital information needed for prediction of the state of node $i$. In contrast, reservoir $R_i$ can compensate for a false positive link from node $j$ to node $i$ by learning, through its training, to ignore its time series input from node $j$. However, if there are too many false positive links to node $i$, reservoir $R_i$ becomes overburdened, and its state prediction accuracy degrades.

*Conclusion.* We are able to construct accurate, data-driven forecasts for the dynamics of large complex networks using a parallel ML architecture that reflects the topology of the network to be predicted. In cases for which a non-parallel approach with comparable resources fails, our scheme is successful when the network links are either known or unknown *a* priori. The parallel nature makes our approach scalable for extremely large networks, creating potential applications to many fields.

---

[1] H. Jaeger and H. Haas, Science **304**, 78 (2004).
[2] Y. Kuramoto, in *International Symposium on Mathematical Problems in Theoretical Physics*, Lecture Notes in Physics (Springer, Berlin, Heidelberg, 1975) pp. 420–422.
[3] J. A. Acebrón, L. L. Bonilla, C. J. Pérez Vicente, F. Ritort, and R. Spigler, Reviews of Modern Physics **77**, 137 (2005).
[4] J. Pathak, B. Hunt, M. Girvan, Z. Lu, and E. Ott, Physical Review Letters **120**, 024102 (2018).
[5] T. Arcomano, I. Szunyogh, J. Pathak, A. Wikner, B. R. Hunt, and E. Ott, Geophysical Research Letters **47**, 10.1029/2020GL087776 (2020).

[6] J. G. Restrepo and E. Ott, EPL (Europhysics Letters) **107**, 60006 (2014).
[7] P. S. Skardal, J. G. Restrepo, and E. Ott, Physical Review E **91**, 060902 (2015).
[8] H. Jaeger, *The 'echo state' approach to analyzing and training recurrent neural networks*, GMD Report 148 (German National Research Center for Information Technology, 2001).
[9] W. Maass, T. Natschläger, and H. Markram, Neural Computation **14**, 2531 (2002).
[10] J. Dambre, D. Verstraeten, B. Schrauwen, and S. Massar, Scientific Reports **2**, 514 (2012).
[11] D. Canaday, A. Griffith, and D. J. Gauthier, Chaos: An Interdisciplinary Journal of Nonlinear Science **28**, 123119 (2018).
[12] P. Holme and M. Huss, Journal of the Royal Society Interface **2**, 327 (2005).
[13] M. Banf and S. Y. Rhee, Scientific Reports **7**, 41174 (2017).
[14] L. Lü and T. Zhou, Physica A: Statistical Mechanics and its Applications **390**, 1150 (2011).
[15] F. Tan, Y. Xia, and B. Zhu, PLOS ONE **9**, e107056 (2014).
[16] S. Kumar and N. Deo, Physical Review E **86**, 026101 (2012).
[17] A. Banerjee, J. Pathak, R. Roy, J. G. Restrepo, and E. Ott, Chaos: An Interdisciplinary Journal of Nonlinear Science **29**, 121104 (2019).
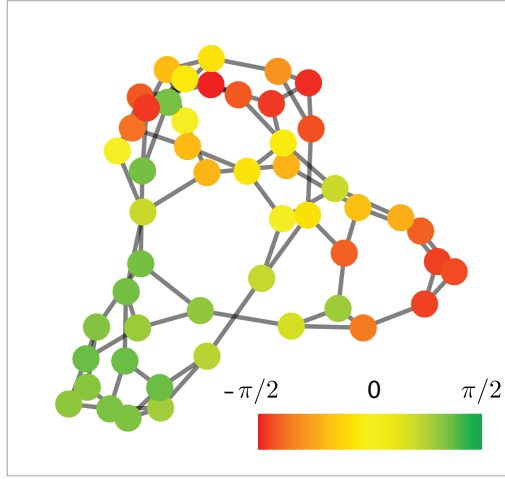[18] T. Schreiber, Physical Review Letters **85**, 461 (2000).

# Supplementary Material

Keshav Srinivasan,[1] Nolan Coble,[1,2] Joy Hamlin,[3] Tom M. Antonsen,[1] Edward Ott,[1] and Michelle Girvan[1]

[1]*University of Maryland, College Park, Maryland 20742, USA.*
[2]*SUNY Brockport, New York 14420, USA.*
[3]*Stony Brook University, New York 11794, USA*

## I.  FREQUENCY ASSORTATIVE KURAMOTO OSCILLATOR NETWORK



Supplementary Figure 1. Frequency assoratative network structure with 50 nodes displaying the natural frequencies of the individual oscillators, with nodes colored by frequency. Similar frequencies are connected with higher probability.
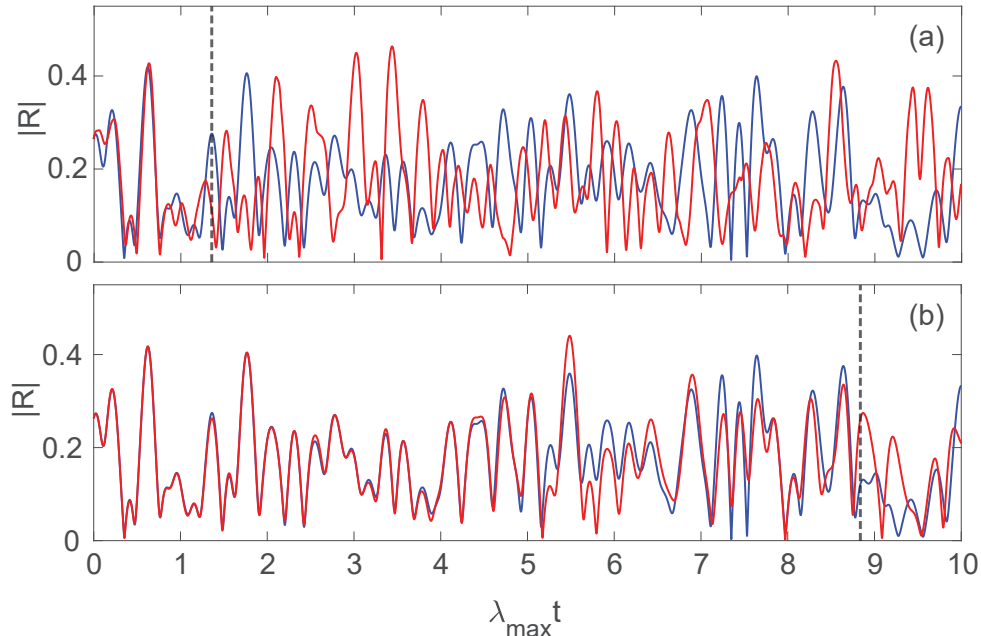
## II.  GENERATING THE RESERVOIR INPUT MATRIX, $\mathbf{W_{in}}$

In this section, we provide more details on how we generated the input matrix, $\mathbf{W_{in}}$, that feeds into the reservoir for predicting our Kuramoto oscillator network. Here the state of each oscillator, $i$, is specified as a tuple $[\sin\theta_i(t), \cos\theta_i(t)]$. The $N_{in}^{osc}$ input oscillators to a particular reservoir are collapsed into a single input vector, $\mathbf{u}$, with $N_{in}^{tot}$ ($N_{in}^{tot} = 2N_{in}^{osc}$) elements such that $\mathbf{u}_i(t) = \sin\theta_i(t)$, $\mathbf{u}_{(N_{in}^{osc}+i)}(t) = \cos\theta_i(t)$ for $1 \le i \le N_{in}^{osc}$. The elements of $\mathbf{W_{in}}$ are chosen so that each of the $N_r$ nodes in the reservoir receives exactly one input from amongst the $2N_{in}^{osc}$ components of the vector $\mathbf{u}(t)$. In the parallel scheme with known neighbors, all the $N_{in}^{tot}$ inputs to a reservoir are treated equally; i.e. each input in $\mathbf{u}(t)$ is connected to the same number ($N_r/N_{in}^{tot}$) of reservoir nodes (Note: $N_r$ is chosen to be an integer multiple of $N_{in}^{tot}$). The non-zero elements of $\mathbf{W_{in}}$ are then randomly drawn from a uniform distribution in $[-\sigma,\sigma]$, where $\sigma$ is the input scaling.

We compare the performance of an all-to-all input matrix to one that is generated as mentioned above, with the input scaling, $\sigma$, optimized separately for each case ($\sigma = 0.3$ for an all-to-all input matrix, for other hyperparameters see Table I). Supp. Fig. 2 clearly highlights the superior performance of the latter method for a Kuramoto oscillator network with all nodes having degree 3, where each reservoir in the parallel scheme (with known neighbors) receives 8 inputs from exactly 4 oscillators: 1 "assigned" + 3 neighbors. The assigned oscillator is the oscillator the reservoir is trained to predict. As we increase the number of inputs to a reservoir, it might be beneficial for each input to be connected to a non-exclusive subset of the reservoir. This would mean each input is still connected to a fixed number of reservoir nodes, but we can relax the criteria that each node in the reservoir is connected to exactly one input.

For parallel prediction with unknown neighbors, we must slightly alter our scheme. In this case, each of the input oscillators to a reservoir are not treated equally. Of the $N_{in}^{osc}$ input oscillators, we only know 1 "true" input oscillator-

the assigned oscillator that reservoir is intended to predict. We cannot distinguish the true neighbors from the false positives in the remaining $N_{in}^{osc} - 1$ input oscillators. We can account for this by "reserving" a larger part of the reservoir for the assigned oscillator and dividing the rest of the reservoir equally between all the inferred neighbors. This means, each of the 2 inputs from the assigned oscillator are connected to $N_{in}^{assign}/2$ reservoir nodes, while the remaining $N_{in}^{tot} - 2$ inputs are each connected to $(N_r - N_{in}^{assign})/(N_{in}^{tot} - 2)$ reservoir nodes. Here $N_r - N_{in}^{assign}$ is chosen to be an integer multiple of $N_{in}^{tot} - 2$ and hence the number of reservoir nodes used in this case can only be set approximately equal to one with known neighbors. For example, in our system with known neighbors we set $N_r = 200$. But for the case with unknown neighbors, where we set $N_{assign} = 50$, if we have 7 inferred neighbors for a particular oscillator, the $N_r$ for the corresponding reservoir must be set to 204 to satisfy the criteria.



Supplementary Figure 2. Parallel prediction of Kuramoto system with known neighbors. The blue curves are the data and red curves are predictions. The dotted lines in each plot denote the valid prediction time. (a) Using an all-to-all input matrix, i.e., all inputs map to all reservoir nodes. (b) Using an input matrix where each input is sent to a different disjoint set of reservoir nodes.
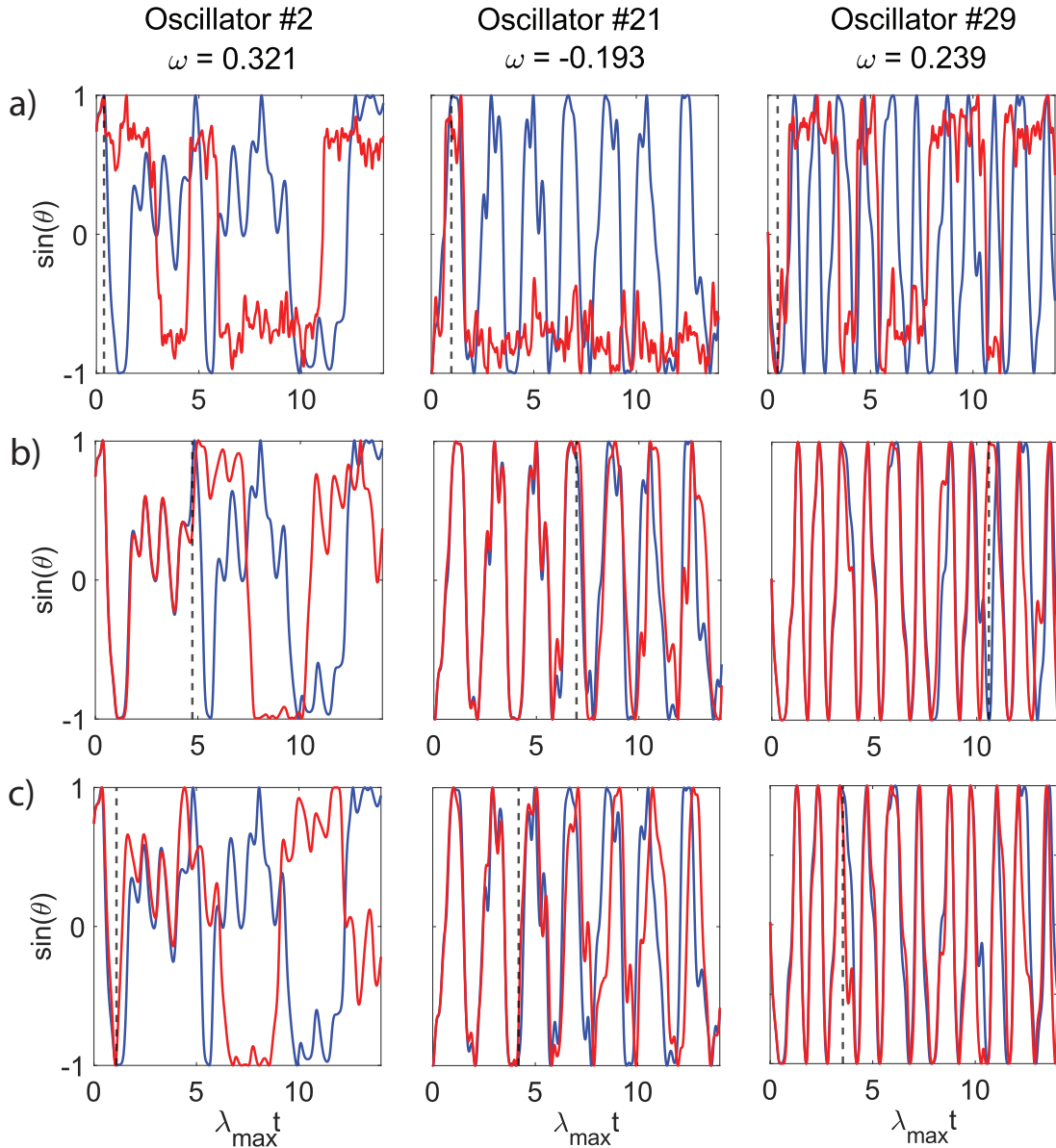
## III.   RC HYPERPARAMETERS

We list the RC hyperparameters (Table I) used to plot the results shown in the paper. These hyperparameters were chosen by a subsequently iterative coarse grid search which approximately maximized the valid prediction time.

| | Number of reservoir nodes | Spectral radius | Input scaling | Leak rate | Regularization parameter |
|---|---|---|---|---|---|
| | $N_r$ | $\rho$ | $\sigma$ | $\alpha$ | $\beta$ |
| Single Reservoir | 10000 | 0.9 | 0.1 | 0 | $10^{-7}$ |
| Parallel Scheme | | | | | |
| a) Known Neighbors | 200 | 0.9 | 0.6 | 0.1 | $10^{-9}$ |
| b) Unknown Neighbors | $\approx 200^*$ | 0.9 | 0.5 | 0.1 | $10^{-8}$ |

TABLE I. Values of reservoir hyperparameters obtained via a coarse grid search for a single reservoir prediction as well as for the parallel reservoir scheme. *Reservoir size set approximately to 200, for details see Section II.

## IV.   NODE-LEVEL PREDICTIONS

While we study the predictive performance of our ML scheme, it is also important to note the node-level details. Supp. Fig. 3(a), shows that a single large reservoir ($N_r = 10000$) is unable to capture the dynamics of individual oscillators in the network. However, our parallel schemes with known edges [Supp. Fig. 3(b)], as well as unknown edges [Supp. Fig. 3(c)], are able to capture the dynamics of the individual oscillators. The parallel schemes, also give good predictions of the system "climate" (i.e. the prediction still seems to resemble the original dynamics) even after the marked valid prediction time (vertical dotted lines in Supp. Fig. 3).



Supplementary Figure 3. Prediction of the dynamics of three representative individual oscillators. The natural frequency of the oscillators is noted above. The blue curves are the data and red curves are predictions. The vertical dotted lines in each plot denote the valid prediction time. (a) Single, non-parallel RC prediction using a large reservoir ($N_r = 10000$), (b) Parallel prediction (with known network links) using 50 separate reservoirs each having modest size ($N_r = 200$), (c) Parallel prediction (with unknown network links) using 50 separate reservoirs each having modest size ($N_r \approx 200$). For row (c) the network structure is estimated by using transfer entropy as a metric for inferring network edges.