

Neural Dynamics on Complex Networks

Chengxi Zang and Fei Wang

Healthcare Policy and Research, Weill Cornell Medicine
chz4001@med.cornell.edu, few2001@med.cornell.edu
New York City, New York, USA

Abstract

Learning dynamics on complex networks governed by differential equation systems is crucial for understanding, predicting, and controlling complex systems in science and engineering. However, this task is very challenging due to the intrinsic complexities in the structures of the high dimensional systems, their elusive continuous-time nonlinear dynamics, and their structural-dynamic dependencies. To address these challenges, we propose a differential deep learning model to learn continuous-time dynamics on complex networks in a data-driven way. We model differential equation systems by graph neural networks. Instead of mapping through a discrete number of hidden layers in the forward process, we solve the initial value problem by integrating the neural differential equation systems over time numerically. In the backward process, we learn the optimal parameters by back-propagating against the forward integration. We validate our model by learning and predicting various real-world dynamics on different complex networks in both (continuous-time) network dynamics learning setting and (regularly-sampled) structured sequence learning setting, and then apply our model to graph semi-supervised classification tasks (a one-snapshot case). The promising experimental results demonstrate our model's capability of jointly capturing the structure, dynamics, and semantics of complex systems in a unified framework.

1 Introduction

Real-world complex systems, such as brain (Gerstner et al. 2014), ecological systems (Gao, Barzel, and Barabási 2016), gene regulation (Alon 2006), human health (Bashan et al. 2016), and social networks (Zang et al. 2018), etc., are usually modeled as complex networks and their evolution are governed by some underlying nonlinear dynamics (Newman, Barabasi, and Watts 2011). Revealing such complex network dynamics is crucial for understanding the complex systems in nature. Effective analytical tools developed for this goal can further help us predict and control these complex systems.

Although the theory of (nonlinear) dynamical systems has been widely studied in different fields including applied math (Strogatz 2018), statistical physics (Newman, Barabasi, and Watts 2011), engineering (Slotine, Li, and others 1991), ecology (Gao, Barzel, and Barabási 2016) and biology (Bashan et al. 2016), these developed models are typically based on a

clear knowledge of the network evolution mechanism which are thus usually referred to as *mechanistic models*. Given the complexity of the real world, there is still a large number of complex networks whose underlying dynamics are unknown yet (e.g., they can be too complex to be modeled by explicit mathematical functions). At the same time, massive data are usually generated during the evolution of these networks. Therefore, modern data-driven approaches are promising and highly demanding in learning the elusive dynamics on complex networks.

The development of a successful data-driven approach for modeling the dynamics on complex networks is very challenging: the interaction structures of the nodes in the network are complex and the number of nodes is large, which is referred to as the high-dimensionality of the complex systems; the rules governing the dynamic change of the nodes' states in complex networks are nonlinear, continuous-time and elusive; the structural and dynamic dependencies within the system are difficult to model. Recently, there is an emerging trend in the data-driven discovery of ordinary differential equations (ODEs) or partial differential equations (PDEs), including sparse regression method (Kutz et al. 2017; Mangan et al. 2016; Rudy et al. 2017), residual network (Qin, Wu, and Xiu 2018), feedforward neural network (Raissi, Perdikaris, and Karniadakis 2018), etc. However, these methods can only handle very small ODE systems or PDEs which consist of only a few terms. Effective learning of the dynamics on large complex networks which consist of tens of thousands of interactions is still largely unknown.

In this paper, we propose a differential deep learning approach to learn continuous-time dynamics on complex networks. We model (high-dimensional) differential equation systems by graph neural networks to capture the instantaneous change of network dynamics. Instead of mapping through a discrete number of layers in the forward process of the conventional neural network models (LeCun, Bengio, and Hinton 2015), we integrate the dynamics on graphs modeled by a neural differential equation system over continuous time. This is like a deep neural network with an infinite number of layers (Chen et al. 2018). In a dynamical system view, the continuous depth can be interpreted as continuous physical time, and the outputs of a hidden layer at time t are instantaneous network dynamics at that moment. In the backward learning process, we back-propagate the gradients of the su-

pervised information w.r.t. the learnable parameters against the forward integration, leading to learning the differential equation system in an end-to-end manner. Besides, we further enhance our algorithm by learning the dynamics in a hidden space learned from the original space of nodes' states. We name our model Neural Dynamics on Complex Networks (NDCN).

We validate our approach by three general tasks: 1) (Network dynamics learning): Can we learn the continuous-time dynamics on complex networks? 2) (Structured sequence learning): Can we predict the regularly-sampled structured sequence? (Seo et al. 2018) 3) (One-snapshot learning): Can we infer the semantic labels of nodes at the terminal time moment? The experimental results show that our model can accurately learn and predict the real-world dynamics on various complex networks, in both continuous-time setting and regularly-sampled sequence setting. What's more, our model learns the semantic labels of nodes in the setting of graph semi-supervised learning (Kipf and Welling 2017) with very competitive performance. Our framework potentially serves as a unified framework to jointly capture the structure, dynamics, and semantics of complex systems in a data-driven manner. Our codes and datasets are open-sourced (Refer to Appendix A).

2 Related work

Dynamics of complex networks. Real-world complex systems are usually modeled as complex networks and driven by nonlinear dynamics: the dynamics of brain and human microbial are examined in (Gerstner et al. 2014) and (Bashan et al. 2016) respectively; (Gao, Barzel, and Barabási 2016) investigated the resilience dynamics of complex systems. (Barzel, Liu, and Barabási 2015) gave a pipeline to construct network dynamics. To the best of our knowledge, our NDCN model is the first neural network approach which learns continuous-time dynamics on complex networks in a data-driven manner.

Data-driven discovery of dynamics. Recently, some data-driven approaches are proposed to learn ODEs or PDEs, including sparse regression (Kutz et al. 2017), residual network (Qin, Wu, and Xiu 2018), feedforward neural network (Raissi, Perdikaris, and Karniadakis 2018), coupled neural networks (Raissi 2018) and so on. (Mangan et al. 2016) tries to learn biological networks dynamics by sparse regression over a large library, which is not scalable to systems with more than 10 nodes. In all, none of them can learn the dynamics on complex systems with more than hundreds of nodes and tens of thousands of interactions.

Neural ODEs. Inspired by residual network (He et al. 2016) and ordinary differential equation (ODE) theory (Lu et al. 2017; Ruthotto and Haber 2018), seminal work neural ODE model (Chen et al. 2018) was proposed to re-write residual networks, normalizing flows, and recurrent neural network in a dynamical system way. However, our NDCN model deals with large and complex differential equations systems. Besides, our model solves different problems, namely learning the dynamics on complex networks.

Optimal control. Relationships between back-propagation in deep learning and optimal control theory are investigated in (Han, Li, and others 2018;

Benning et al. 2019). We formulate our loss function by leveraging the concept of running loss and terminal loss in optimal control. We give novel constraints in optimal control which is modeled by neural differential equations systems on graphs. Our model solves novel tasks, e.g. learning the dynamics on complex networks and refer to Sec.3.1.

Graph neural networks and temporal-GNNs. Graph neural networks (GNNs) (Wu et al. 2019b), e.g., Graph convolution network (GCN) (Kipf and Welling 2017), attention-based GNN (AGNN) (Thekumparampil et al. 2018), graph attention networks (GAT) (Veličković et al. 2017), etc., achieved state-of-the-art performance on graph semi-supervised learning tasks. However, existing GNNs usually have 1 or 2 layers and can not go deep (Li, Han, and Wu 2018; Wu et al. 2019b). Our NDCN gives a dynamical system view on GNNs: the continuous depth can be interpreted as continuous physical time, and the outputs of a hidden layer at time t are instantaneous network dynamics at that moment. By capturing continuous-time network dynamics and their transient behaviors, our model gives very competitive and even better results than above GNNs.

By combining RNNs or convolution operators with GNNs, temporal-GNNs (Yu, Yin, and Zhu 2017; Kazemi et al. 2019; Narayan and Roe 2018; Seo et al. 2018) try to predict the regularly-sampled structured sequences. However, these models can not be applied to continuous-time dynamics (observed at arbitrary physical times with different time intervals). Our NDCN not only predicts the continuous-time network dynamics at an arbitrary time or semantic labels from one snapshot but also predicts the structured sequences very well in a more succinct way with much fewer parameters.

3 Preliminaries

3.1 Problem Definition

We first introduce a differential equation system which models the dynamics on complex networks:

$$\frac{dX(t)}{dt} = f(X(t); G; W(t); t); \quad (1)$$

where $X(t) \in \mathbb{R}^{n \times d}$ represents the state (node feature values) of a dynamic system consisting of n linked nodes at time $t \in [0, 1)$, and each node is characterized by d dimensional features. $G = (V, E)$ is the network structure capturing how the nodes are linked to each other. $W(t)$ are parameters which control how the system evolves over time. $X(0) = X_0$ is the initial states of this system at time $t = 0$. The function $f : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ is a function governing the dynamics of the system, which could be either linear or nonlinear. In addition, nodes can have various semantic labels $Y(X, t) \in \{0, 1\}^k$ at time t , and θ represents the parameters of this classification function. The problems we are trying to solve in this paper are:

(Network dynamics learning) How to learn the continuous-time dynamics $\frac{dX(t)}{dt}$ on complex networks from empirical data? Mathematically, given a graph G and the observations of the states of system $\{X(\hat{t}_1), X(\hat{t}_2), \dots, X(\hat{t}_T)\}_{0 \leq \hat{t}_1 < \dots < \hat{t}_T \leq 1}$, and t_1 to t_T are arbitrary physical time stamps, can we learn

differential equation systems $\frac{dX(t)}{dt} = f(X, G, W, t)$ to generate or predict continuous-time dynamics $X(t)$ at arbitrary physical time t ? The arbitrary physical time means that t_1, \dots, t_T are possibly irregularly sampled with different observational time intervals. When $t > t_T$, we call the task **extrapolation prediction**, while $t < t_T$ and $t \notin \{t_1, \dots, t_T\}$ for **interpolation prediction**.

(Structured sequence learning). As a special case when t_1, t_2, \dots, t_T are sampled regularly with equal time intervals, the above problem degenerates to a structured sequence learning task with an emphasis on sequential order instead of arbitrary physical time. The goal is to extrapolate next m steps' $X[t_T + 1], \dots, X[t_T + m]$.

(One-snapshot learning) How to learn the semantic labels of $Y(X(t_T))$ at the moment $t = t_T$ for each node?

As a special case of above problem with an emphasis on a specific moment and without loss of generality, we focus on the moment at the terminal time t_T . The function Y can be a mapping from the nodes' states (e.g. humidity) to their labels (e.g. taking umbrella or not).

3.2 Network Dynamics

We investigate the following three real-world network dynamics. Let $x_i(t) \in \mathbb{R}^d$ be d dimensional features of node i at time t and thus $X(t) = [\dots, x_i(t), \dots]^T \in \mathbb{R}^{n \times d}$. We show their differential equation systems in vector form for clarity and implement them in matrix form:

The heat diffusion dynamics $\frac{dx_i(t)}{dt} = k_{ij} \sum_{j=1}^n A_{ij} (x_j - x_i)$ governed by Newton's law of cooling (Luikov 2012), which states that the rate of heat change of node i is proportional to the difference of the temperature between node i and its neighbors with heat capacity matrix A .

The mutualistic interaction dynamics among species in ecology, governed by equation $\frac{dx_i(t)}{dt} = b_i + \frac{x_i}{c_i} (1 - \sum_{j=1}^n A_{ij} \frac{x_j}{d_i + e_i x_i + h_j x_j})$ (For brevity, the operations between vectors are element-wise). The mutualistic differential equation systems (Gao, Barzel, and Barabási 2016) capture the abundance $x_i(t)$ of species i , consisting of incoming migration term b_i , logistic growth with population capacity k_i (Zang et al. 2018) and Allee effect (Allee et al. 1949) with cold-start threshold c_i , and mutualistic interaction term with interaction network A .

The gene regulatory dynamics governed by Michaelis-Menten equation $\frac{dx_i(t)}{dt} = b_i x_i^f + \sum_{j=1}^n A_{ij} \frac{x_j^h}{x_j^h + 1}$ where the first term models degradation when $f = 1$ or dimerization when $f = 2$, and the second term captures genetic activation tuned by the Hill coefficient h (Alon 2006; Gao, Barzel, and Barabási 2016).

Complex Networks. We consider following networks: (a) Grid network, where each node is connected with 8 neighbors (as shown in Fig. 2(a)); (b) Random network, generated by Erdős and Rényi model (Erdos and Renyi 1959) (as shown in Fig. 2(b)); (c) Power-law network, generated by Albert-Barabási model (Barabási and Albert 1999) (as

shown in Fig. 2(c)); (d) Small-world network, generated by Watts-Strogatz model (Watts and Strogatz 1998) (as shown in Fig. 2(d)); and (e) Community network, generated by random partition model (Fortunato 2010) (as shown in Fig. 2(e)).

Visualization. To visualize dynamics on complex networks over time is not trivial. We first generate a network with n nodes by aforementioned network models. The nodes are re-ordered according to the community detection method by Newman (Newman 2010) and each node has a unique label from 1 to n . We layout these nodes on a 2-dimensional $\frac{n}{n} \times \frac{n}{n}$ grid and each grid point $(r, c) \in \mathbb{N}^2$ represents the i^{th} node where $i = r \frac{n}{n} + c + 1$. Thus, nodes' states $X(t) \in \mathbb{R}^{n \times d}$ at time t when $d = 1$ can be visualized as a scalar field function $X : \mathbb{N}^2 \rightarrow \mathbb{R}$ over the grid. Please refer to Appendix B for the animations of these dynamics on different complex networks over time.

4 General framework

We formulate our general framework as follows:

$$\begin{aligned} \underset{W(t)}{\operatorname{argmin}} \quad & L = \int_0^T R(X(t); G; W; t) dt + S(Y(X(T), \cdot)) \\ \text{subject to} \quad & \frac{dX(t)}{dt} = f(X(t); G; W; t); X_0 \end{aligned} \quad (2)$$

where $R(X(t), G, W, t)$ is the running loss of the dynamics on graph at time t , and $S(Y(X(T), \cdot))$ is the terminal semantic loss at time T . By integrating $\frac{dX}{dt} = f(X, G, W, t)$ over time t from initial state X_0 , a.k.a. solving the initial value problem (Boyce, DiPrima, and Meade 1992) for this differential equation system, we can get the continuous-time dynamics $X(t) = X(0) + \int_0^t f(X(\tau), G, W, \tau) d\tau$ at arbitrary time moment $t > 0$.

Such a formulation can be seen as an optimal control problem so that the goal becomes to obtain the best control parameters $W(t)$ for differential equation system $\frac{dX}{dt} = f(X, G, W, t)$ and the best classification parameters for semantic function $Y(X(t), \cdot)$ by solving above optimization problem. Different from traditional optimal control framework, we model the differential equation systems $\frac{dX}{dt} = f(X, G, W, t)$ by graph neural networks. By integrating $\frac{dX}{dt} = f(X, G, W, t)$ over continuous time, namely $X(t) = X(0) + \int_0^t f(X(\tau), G, W, \tau) d\tau$, we get our differential deep learning models. In a dynamical system view, our differential deep learning models can be a time-varying coefficient dynamical system when $W(t)$ changes over time; or a constant coefficient dynamical system when W is constant over time for parameter sharing. It's worthwhile to recall that the deep learning methods with L hidden neural layers f are $X[L] = f_L \dots f_2 f_1(X[0])$, which are iterated maps (Strogatz 2018) with an integer number of discrete layers and thus can not learn continuous-time dynamics $X(t)$ at arbitrary time. In contrast, our model $X(t) = X(0) + \int_0^t f(X(\tau), G, W, \tau) d\tau$ can have continuous layers with a real number t depth corresponding to continuous-time dynamics.

Moreover, to further increase the express ability of our model, we can encode the network signal $X(t)$ from the original space to $X_h(t)$ in hidden space (usually with a different

number of dimensions), and learn the dynamics in such a space. Then our model becomes:

$$\begin{aligned} \operatorname{argmin}_{W(t); (T)} \quad & L = \int_0^T R(X(t); G; W; t) dt + S(Y(X(T));) \\ \text{subject to} \quad & X_h(t) = f_e(X(t)) \\ & \frac{dX_h(t)}{dt} = f(X_h(t); G; W; t; X_h(0)) \\ & X(t) = f_d(X_h(t)) \end{aligned} \quad (3)$$

where the first constraint transforms $X(t)$ into hidden space $X_h(t)$ through encoding function f_e . The second constraint is the governing dynamics in the hidden space. The third constraint decodes the hidden signal back to the original space with decoding function f_d . The design of f_e , f , and f_d are flexible to be any neural structure (e.g. f_d is a softmax function for classification). We denote our model as *Neural Dynamics on Complex Networks (NDCN)*.

We solve the initial value problem (i.e., integrating the differential equation systems over time numerically) by numerical methods (e.g., 1st-order Euler method, high-order method Dormand-Prince DOPRI5 (Dormand 1996), etc.). The numerical methods can approximate continuous-time dynamics $X(t) = X(0) + \int_0^t f(X(\tau), G, W, \tau) d\tau$ at arbitrary time t accurately with guaranteed error. In order to learn the learnable parameters W , we back-propagate the gradients of the loss function w.r.t the control parameters $\frac{\partial L}{\partial W}$ over the numerical integration process backwards in an end-to-end manner, and solve the optimization problem by stochastic gradient descent methods (e.g., Adam (Kingma and Ba 2015)). We will show some concrete examples of the above framework in the next three sections.

5 Learning continuous-time network dynamics

In this section, we investigate if our NDCN model can learn aforementioned network dynamics in the continuous-time setting for both interpolation prediction and extrapolation prediction. The continuous-time setting means that the observational times t_1 to t_T of the observed states of system $fX(\hat{t}_1), X(\hat{t}_2), \dots, X(\hat{t}_T)g$ are arbitrary physical time stamps which are irregularly sampled with different observational time intervals. The extrapolation prediction is to predict $X(t)$ at arbitrary physical time moment t when $t > t_T$, while the interpolation prediction is to predict $X(t)$ when $t < t_T$ and $t \notin \{t_1, \dots, t_T\}$.

5.1 A Model Instance

We solve the objective function in (3) with an emphasis on running loss only. Without the loss of generality, we use ℓ_1 -norm loss as the running loss R . More concretely, we adopt two fully connected neural layers with a nonlinear hidden layer as the encoding function f_e , a graph convolution neural network (GCN) like structure (Kipf and Welling 2017) but with a different graph diffusion operator to model the instantaneous network dynamics in the hidden space, and a linear decoding function f_d for regression tasks in the

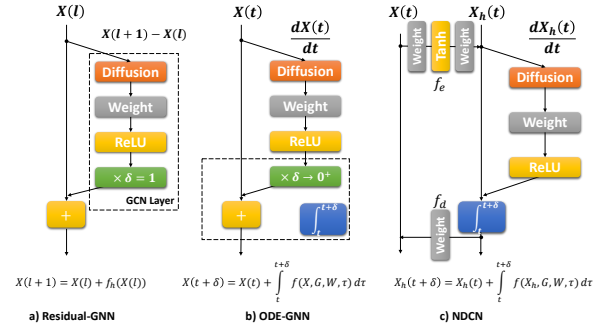


Figure 1: *Illustration of an NDCN instance.* a) Residual Graph Neural Networks, b) ODE-GNN model and c) Our Neural Dynamics on Complex Network (NDCN) model. The integer l represents the discrete l^{th} layer and the real number t represents continuous physical time.

original signal space. Thus, our model is (see neural structure in Figure 1c):

$$\begin{aligned} \operatorname{argmin}_{W_e, b_e} \quad & L = \int_0^T jX(t) \hat{X}(t) dt \\ \text{subject to} \quad & X_h(t) = \tanh(X(t)W_e + b_e)W_0 + b_0 \\ & \frac{dX_h(t)}{dt} = \text{ReLU}(X_h(t)W + b) ; X_h(0) \\ & X(t) = X_h(t)W_d + b_d \end{aligned} \quad (4)$$

where $\hat{X}(t) \in \mathbb{R}^n$ is the supervised dynamic information available at time stamp t (in the semi-supervised case the missing information can be padded by 0). The j denotes ℓ_1 -norm loss (mean element-wise absolute value difference) between $X(t)$ and $\hat{X}(t)$ at time $t \in [0, T]$. We adopt diffusion operator $\Delta = D^{-\frac{1}{2}}(D - A)D^{-\frac{1}{2}} \in \mathbb{R}^{n \times n}$ which is the normalized graph Laplacian where $A \in \mathbb{R}^{n \times n}$ is the adjacency matrix of the network and $D \in \mathbb{R}^{n \times n}$ is the corresponding node degree matrix. The $W \in \mathbb{R}^{d_e \times d_e}$ and $b \in \mathbb{R}^{d_e}$ are shared parameters (namely, the weights and bias of a linear connection layer) over time $t \in [0, T]$. The $W_e \in \mathbb{R}^{d \times d_e}$ and $W_0 \in \mathbb{R}^{d_e \times d_e}$ are the matrices in linear layers for encoding, while $W_d \in \mathbb{R}^{d_e \times d}$ are for decoding. The b_e, b_0, b, b_d are the biases at the corresponding layer. We learn the parameters $W_e, W_0, W, W_d, b_e, b_0, b, b_d$ from empirical data so that we can learn X in a data-driven manner.

We design the neural differential equation system as $\frac{dX(t)}{dt} = \text{ReLU}(X(t)W + b)$ to learn any unknown network dynamics. We can regard $\frac{dX(t)}{dt}$ as a single neural layer at time moment t . The $X(t)$ at arbitrary time t is achieved by integrating $\frac{dX(\tau)}{d\tau}$ over time, i.e., $X(t) = X(0) + \int_0^t \text{ReLU}(X(\tau)W + b) d\tau$, leading to a continuous-time deep neural network.

5.2 Experiments

Baselines. To the best of our knowledge, there are no baselines for learning continuous-time dynamics on complex networks, and thus we compare the ablation models of NDCN for this task. By investigating ablation models we show that our NDCN is a minimum model for this task. We keep the loss function the same and construct the following baselines:

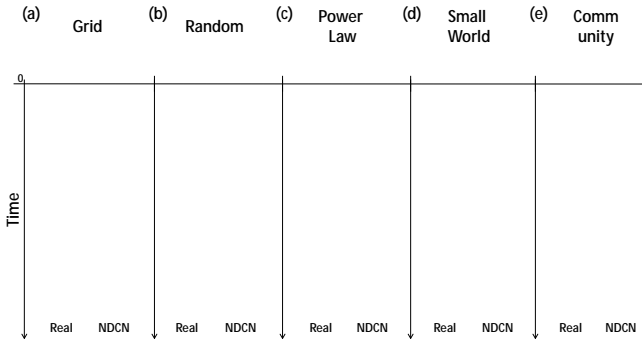


Figure 2: *Heat diffusion on different networks.* Each of the five vertical panels represents the dynamics on one network over physical time. For each network dynamics, we illustrate the sampled ground truth dynamics (left) and the dynamics generated by our NDCN (right) from top to bottom following the direction of time.

The model without encoding f_e and f_d and thus no hidden space: $\frac{dX(t)}{dt} = \text{ReLU}(X(t)W + b)$, namely ODE-GNN, which learns the dynamics in the original signal space $X(t)$ as shown in Fig. 1b;

The model without graph diffusion operator: $\frac{dX_h(t)}{dt} = \text{ReLU}(X_h(t)W + b)$, i.e., an ODE Neural Network, which can be thought as a continuous-time version of forward residual neural network (See Fig. 1a and Fig. 1b for the difference between residual network and ODE network).

The model without control parameters W : $\frac{dX_h(t)}{dt} = \text{ReLU}(X_h(t))$ which has no linear connection layer between t and $t + dt$ (where $dt \neq 0$) and thus indicating a determined dynamics to spread signals.

Experimental setup. We generate underlying networks with 400 nodes by network models in Sec.3.2 and the illustrations are shown in Fig. 2,3 and 4. We set the initial value $X(0)$ the same for all the experiments and thus different dynamics are only due to their different dynamic rules and underlying networks (See Appendix B).

We irregularly sample 120 snapshots of the continuous-time dynamics $fX(\hat{t}_1), \dots, X(\hat{t}_{120})/0 < t_1 < \dots < t_{120}$ Tg where the time intervals between t_1, \dots, t_{120} are different. We randomly choose 80 snapshots from $X(\hat{t}_1)$ to $X(\hat{t}_{100})$ for training, the left 20 snapshots from $X(\hat{t}_1)$ to $X(\hat{t}_{100})$ for testing the interpolation prediction task. We use the 20 snapshots from $X(\hat{t}_{101})$ to $X(\hat{t}_{120})$ for testing the extrapolation prediction task.

We use Dormand-Prince method (Dormand 1996) to get the ground truth dynamics, and use Euler method in the forward process of our NDCN (More configurations in Appendix C). We evaluate the results by ℓ_1 loss and normalized ℓ_1 loss (normalized by the mean element-wise value of $X(\hat{t})$), and they lead to the same conclusion (We report normalized ℓ_1 loss here and see Appendix E for ℓ_1 loss). Results are the mean and standard deviation of the loss over 20 independent runs for 3 dynamic laws on 5 different networks by each method.

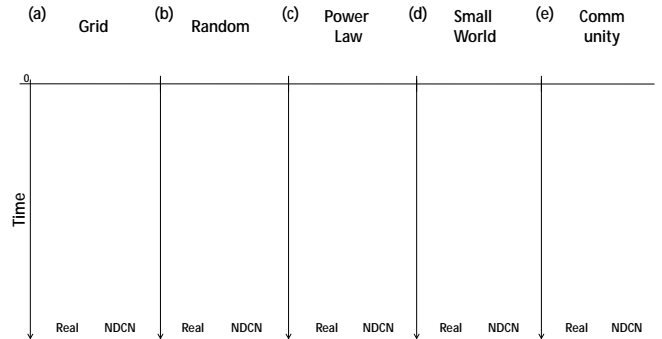


Figure 3: *Biological mutualistic interaction on different networks.*

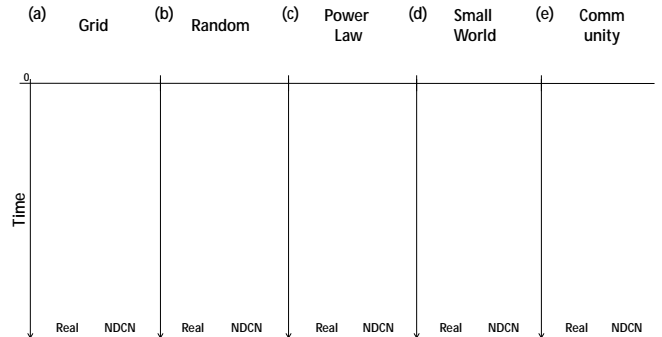


Figure 4: *Gene regulation dynamics on different networks.*

Results We visualize the ground-truth and learned dynamics in Fig. 2,3 and 4, and please see the animations of these network dynamics in Appendix B. We find that one dynamic law may behave quite different on different networks: heat dynamics may gradually die out to be stable but follow different dynamic patterns in Fig. 2. Gene dynamics are asymptotically stable on grid in Fig. 4a but unstable on random networks in Fig. 4b or community networks in Fig. 4e. Both gene regulation dynamics in Fig. 4c and biological mutualistic dynamics in Fig. 3c show very bursty patterns on power-law networks. However, visually speaking, our NDCN learns all these different network dynamics very well.

The quantitative results of extrapolation and interpolation prediction are summarized in Table 1 and Table 2 respectively. We observe that our NDCN captures different dynamics on various complex networks accurately and outperforms all the continuous-time baselines by a large margin, indicating that our NDCN potentially serves as a minimum model in learning continuous-time dynamics on complex networks.

6 Learning Regularly-sampled Dynamics

What's more, our model also captures dynamics from regularly-sampled network dynamics (i.e. the structured sequence learning setting) very well.

Baselines. We compare our model with the temporal-GNN models which are usually combinations of RNN models and GNN models (Kazemi et al. 2019; Narayan and Roe 2018; Seo et al. 2018). We use GCN (Kipf and Welling 2017) as a graph structure extractor and use LSTM/GRU/RNN (Lipton,

Table 1: **Continuous-time Extrapolation Prediction.** Our NDCN predicts different continuous-time network dynamics accurately. Each result is the normalized ℓ_1 error with standard deviation (in percentage %) from 20 runs for 3 dynamics on 5 networks by each method.

		Grid		Random		Power Law		Small World		Community	
Heat Diffusion	No-Encode	29:9	7:3	27:8	5:1	24:9	5:2	24:8	3:2	30:2	4:4
	No-Graph	30:5	1:7	5:8	1:3	6:8	0:5	10:7	0:6	24:3	3:0
	No-Control	73:4	14:4	28:2	4:0	25:2	4:3	30:8	4:7	37:1	3:7
	NDCN	4:1	1:2	4:3	1:6	4:9	0:5	2:5	0:4	4:8	1:0
Mutualistic Interaction	No-Encode	45:3	3:7	9:1	2:9	29:9	8:8	54:5	3:6	14:5	5:0
	No-Graph	56:4	1:1	6:7	2:8	14:8	6:3	54:5	1:0	9:5	1:5
	No-Control	140:7	13:0	10:8	4:3	106:2	42:6	115:8	12:9	16:9	3:1
	NDCN	26:7	4:7	3:8	1:8	7:4	2:6	14:4	3:3	3:6	1:5
Gene Regulation	No-Encode	31:7	14:1	17:5	13:0	33:7	9:9	25:5	7:0	26:3	10:4
	No-Graph	13:3	0:9	12:2	0:2	43:7	0:3	15:4	0:3	19:6	0:5
	No-Control	65:2	14:2	68:2	6:6	70:3	7:7	58:6	17:4	64:2	7:0
	NDCN	16:0	7:2	1:8	0:5	3:6	0:9	4:3	0:9	2:5	0:6

Table 2: **Continuous-time Interpolation Prediction.** Our NDCN predicts different continuous-time network dynamics accurately. Each result is the normalized ℓ_1 error with standard deviation (in percentage %) from 20 runs for 3 dynamics on 5 networks by each method.

		Grid		Random		Power Law		Small World		Community	
Heat Diffusion	No-Encode	32:0	12:7	26:7	4:4	25:7	3:8	27:9	7:3	35:0	6:3
	No-Graph	41:9	1:8	9:4	0:6	18:2	1:5	25:0	2:1	25:0	1:4
	No-Control	56:8	2:8	32:2	7:0	33:5	5:7	40:4	3:4	39:1	4:5
	NDCN	3:2	0:6	3:2	0:4	5:6	0:6	3:4	0:4	4:3	0:5
Mutualistic Interaction	No-Encode	28:9	2:0	19:9	6:5	34:5	13:4	27:6	2:6	25:5	8:7
	No-Graph	28:7	4:5	7:8	2:4	23:2	4:2	26:9	3:8	14:1	2:4
	No-Control	72:2	4:1	22:5	10:2	63:8	3:9	67:9	2:9	33:9	12:3
	NDCN	7:6	1:1	6:6	2:4	6:5	1:3	4:7	0:7	7:9	2:9
Gene Regulation	No-Encode	39:2	13:0	14:5	12:4	33:6	10:1	27:7	9:4	21:2	10:4
	No-Graph	25:2	2:3	11:9	0:2	39:4	1:3	15:7	0:7	18:9	0:3
	No-Control	66:9	8:8	31:7	5:2	40:3	6:6	49:0	8:0	35:5	5:3
	NDCN	5:8	1:0	1:5	0:6	2:9	0:5	4:2	0:9	2:3	0:6

Berkowitz, and Elkan 2015) to learn the temporal relationships between ordered structured sequences. We keep the loss function the same and construct the following baselines. We denote each recurrent cell as LSTM/GRU/RNN and refer to Appendix D for the detailed equations.

LSTM-GNN: the temporal-GNN with LSTM cell $X[t + 1] = LSTM(GCN(X[t], G))$.

GRU-GNN: the temporal-GNN with GRU cell $X[t + 1] = GRU(GCN(X[t], G))$.

RNN-GNN: the temporal-GNN with RNN cell $X[t + 1] = RNN(GCN(X[t], G))$.

Experimental setup. We regularly sample 100 snapshots of the continuous-time network dynamics $\{X[\hat{t}_1], \dots, X[\hat{t}_{100}]\}_{0 \leq t_1 < \dots < t_{120} \leq T}$ where the time intervals between t_1, \dots, t_{100} are the same. We use first 80 snapshots $X[\hat{t}_1], \dots, X[\hat{t}_{80}]$ for training and the left 20 snapshots $X[\hat{t}_{81}], \dots, X[\hat{t}_{100}]$ for testing extrapolation prediction task. The temporal-GNN models are usually used for next few step prediction and cannot be used for the interpolation task (say, to predict $X[t_{1.23}]$) directly. We use 5 and 10 for hidden dimension of GCN and RNN models respectively. Other settings are the same as previous continuous-time dynamics experiment.

Results We summarize the results of the extrapolation prediction of regularly-sampled dynamics in Table 3. The GRU-GNN model works well in mutualistic dynamics on random network and community network. Our NDCN predicts different dynamics on these complex networks accurately and outperforms the baselines in almost all the settings. What’s more, our model capture the structure and dynamics in a much more succinct way. The learnable parameters of our NDCN, RNN-GNN, GRU-GNN, LSTM-GNN are **901**, 24530, 64770, and

84890 respectively.

7 Learning semantic labels at terminal time

We investigate the third question, i.e., how to learn the semantic labels of each node at the terminal time? Various graph neural networks (GNN) (Wu et al. 2019b) achieve the state-of-the-art performance in graph semi-supervised classification task (Yang, Cohen, and Salakhutdinov 2016; Kipf and Welling 2017). Existing GNNs usually adopt 1 or 2 hidden layers (Kipf and Welling 2017; Veličković et al. 2017) and cannot go deep (Li, Han, and Wu 2018). Our framework follows the perspective of a dynamical system, and goes beyond an integer number L of hidden layers in GNNs to a real number depth t of hidden layers, implying continuous-time dynamics on the graph. By integrating continuous-time dynamics on the graph over time, we get a more fine-grained forward process and thus our NDCN model shows very competitive even better results compared with state-of-the-art GNN models which may have sophisticated parameters (e.g. attention).

7.1 A Model Instance

Following the same framework as in Section 3, we propose a simple model with the terminal semantic loss $S(Y(T))$ modeled by the cross-entropy loss for classification task:

$$\begin{aligned} \arg\min_{W_e, b_e; W_d, b_d} L &= \int_0^T R(t) dt \quad \sum_{i=1}^n \sum_{k=1}^c \hat{Y}_{i:k}(T) \log Y_{i:k}(T) \\ \text{subject to} \quad X_h(0) &= \tanh(X(0)W_e + b_e) \\ \frac{dX_h(t)}{dt} &= \text{ReLU}(X_h(t)) \\ Y(T) &= \text{softmax}(X_h(T)W_d + b_d) \end{aligned} \quad (5)$$

where $Y(T) \in \mathbb{R}^{n \times c}$ is the label distributions of nodes at time $T \in \mathbb{R}$ whose element $\hat{Y}_{i:k}(T)$ denotes the probability

Table 3: Regularly-sampled Extrapolation Prediction. Our NDCN predicts different structured sequences accurately. Each result is the normalized error with standard deviation (in percentage) from 20 runs for 3 dynamics on 5 networks by each method.

		Grid		Random		Power Law		Small World		Community	
Heat Diffusion	LSTM-GNN	12:8	2:1	21:6	7:7	12:4	5:1	11:6	2:2	13:5	4:2
	GRU-GNN	11:2	2:2	9:1	2:3	8:8	1:3	9:3	1:7	7:9	0:8
	RNN-GNN	18:8	5:9	25:0	5:6	18:9	6:5	21:8	3:8	16:1	0:0
	NDCN	4:3	0:7	4:7	1:7	5:4	0:4	2:7	0:4	5:3	0:7
Mutualistic Interaction	LSTM-GNN	51:4	3:3	24:2	24:2	27:0	7:1	58:2	2:4	25:0	22:3
	GRU-GNN	49:8	4:1	1:0	3:6	12:2	0:8	51:1	4:7	3:7	4:0
	RNN-GNN	56:6	0:1	8:4	11:3	12:0	0:4	57:4	1:9	8:2	6:4
	NDCN	29:8	1:6	4:7	1:1	11:2	5:0	15:9	2:2	3:8	0:9
Gene Regulation	LSTM-GNN	27:7	3:2	67:3	14:2	38:8	12:7	13:1	2:0	53:1	16:4
	GRU-GNN	24:2	2:8	50:9	6:4	35:1	15:1	11:1	1:8	46:2	7:6
	RNN-GNN	28:0	6:8	56:5	5:7	42:0	12:8	14:0	5:3	46:5	3:5
	NDCN	18:6	9:9	2:4	0:9	4:1	1:4	5:5	0:8	2:9	0:5

of the nodes $i = 1; \dots; n$ with labels $k = 1; \dots; c$ at time T . The $\hat{Y} \in \mathbb{R}^{n \times c}$ is the supervised information (again missing information can be padded by 0) observed at time T . We use differential equation system $\frac{dx(t)}{dt} = \text{ReLU}(X(t))$ to spread the graph signals over continuous time, i.e., $X_h(T) = X_h(0) + \int_0^T \text{ReLU}(X_h(t)) dt$.

Compared with the model in Eq.4, we only have supervised information from one snapshot at time T . Thus, we model the running loss $\int_0^T R(t) dt$ as the L_2 -norm regularizer of the learnable parameters $\int_0^T R(t) dt = (jW_e j_2^2 + jW_d j_2^2 + jW_a j_2^2 + jW_b j_2^2)$ to avoid overfitting. We adopt the diffusion operator $\mathcal{D} = \frac{1}{2}(I + (1 - \alpha)A)\mathcal{D}^{\frac{1}{2}}$ where A is the adjacency matrix, D is the degree matrix and $\mathcal{D} = I + (1 - \alpha)D$ keeps \mathcal{D} normalized. The parameter $\alpha \in [0, 1]$ tunes nodes' adherence to their previous information or their neighbors' collective opinion. We use it as a hyper-parameter here for simplicity and we can make it as a learnable parameter later. The differential equation system $\frac{dx}{dt} = X$ follows the dynamics of averaging the neighborhood opinion as $\frac{dx_i(t)}{dt} = \frac{1}{(1 - \alpha)d_i + 1} x_i(t) + \frac{\alpha}{(1 - \alpha)d_i + 1} \sum_j A_{ij} \frac{1}{(1 - \alpha)d_j + 1} x_j(t)$ for node i . When $\alpha = 0$, $\frac{dx_i(t)}{dt} = x_i(t)$ averages the neighbors as normalized random walk, when $\alpha = 1$, $\frac{dx_i(t)}{dt} = x_i(t)$ captures exponential dynamics without network effects, and when $\alpha = 0.5$, $\frac{dx_i(t)}{dt} = x_i(t)$ averages both neighbors and itself as in (Kipf and Welling 2017).

7.2 Experiments

We validate our model in the graph semi-supervised classification setting. For the consistency of comparison with prior works, we follow the same experimental setup as (Kipf and Welling 2017; Velicković et al 2017; Thekumparampil et al. 2018). Refer to Appendix F for the detailed information about datasets, baselines, and their configurations.

Results We summarize the results in Table 4. We find our NDCN outperforms many state-of-the-art GNN models. Results for the baselines are taken from (Kipf and Welling 2017; Velicković et al 2017; Thekumparampil et al. 2018; Wu et al. 2019a). We report the mean and standard deviation of our results for 100 runs. We get our reported results in Table 4 when terminal time $T = 1:2$, $\alpha = 0$ for the Cora dataset, $T = 1:0$, $\alpha = 0:8$ for the Citeseer dataset, and $T = 1:1$, $\alpha = 0:4$ for the Pubmed dataset.

Table 4: Test mean accuracy with standard deviation in percentage (%) over 100 runs. Our NDCN model gives very competitive results compared with many GNN models.

Model	Cora		Citeseer		Pubmed	
GCN	81:5		70:3		79:0	
AGNN	83:1	0:1	71:7	0:1	79:9	0:1
GAT	83:0	0:7	72:5	0:7	79:0	0:3
NDCN	83:3	0:6	73:1	0:6	79:8	0:4

Figure 5: Our NDCN model captures continuous-time dynamics. Mean classification accuracy of 100 runs over terminal time when given a specific T . Insets are the accuracy over the two-dimensional space of terminal time and α .

By capturing the continuous-time network dynamics to diffuse network signals, our NDCN gives better classification accuracy at terminal time $T \in \mathbb{R}^+$. Figure 5 plots the mean accuracy with error bars over terminal time in the abovementioned settings (we further plot the accuracy over terminal time T and α in the insets and Appendix G). We find for all the three datasets their accuracy curves follow rise and fall patterns around the best terminal time. Indeed, when the terminal time T is too small or too large, the accuracy degenerates because the features of nodes are in under-diffusion or over-diffusion states, implying the necessity in capturing continuous-time dynamics. In contrast, previous GNNs can only have a discrete number of layers which can not capture the continuous-time network dynamics accurately.

8 Conclusion

We propose a differential deep learning model to learn continuous-time dynamics on complex networks. We model differential equations systems by graph neural networks and integrate the neural differential equations systems over time. By capturing the continuous-time network dynamics, our NDCN gives the meanings of physical time and the continuous-time network dynamics to the depth and hidden outputs respectively, learns real-world dynamics on complex network accurately in both (irregularly-sampled) continuous-time setting and (regularly-sampled) structured sequence setting, and outperforms many GNN models in the graph semi-supervised classification task (a one-snapshot case). Codes and datasets are open-sourced (See Appendix A).

References

- Allee, W. C.; Park, O.; Emerson, A. E.; Park, T.; Schmidt, K. P.; et al. 1949. Principles of animal ecology. Technical report, Saunders Company Philadelphia, Pennsylvania, USA.
- Alon, U. 2006. An introduction to systems biology: design principles of biological circuits. Chapman and Hall/CRC.
- Barabási, A.-L., and Albert, R. 1999. Emergence of scaling in random networks. *Science* 286(5439):509–512.
- Barzel, B.; Liu, Y.-Y.; and Barabási, A.-L. 2015. Constructing minimal models for complex system dynamics. *Nature communications* 6:1–10.
- Bashan, A.; Gibson, T. E.; Friedman, J.; Carey, V. J.; Weiss, S. T.; Hohmann, E. L.; and Liu, Y.-Y. 2016. Universality of human microbial dynamics. *Nature* 534(7606):259.
- Benning, M.; Celledoni, E.; Ehrhardt, M. J.; Owren, B.; and Schönlieb, C.-B. 2019. Deep learning as optimal control problems: models and numerical methods. *arXiv preprint arXiv:1904.05657*
- Boyce, W. E.; DiPrima, R. C.; and Meade, D. B. 1992. Elementary differential equations and boundary value problems. volume 9. Wiley New York.
- Chen, T. Q.; Rubanova, Y.; Bettencourt, J.; and Duvenaud, D. K. 2018. Neural ordinary differential equations. *Advances in Neural Information Processing Systems* 31:6571–6583.
- Dormand, J. R. 1996. Numerical methods for differential equations: a computational approach. volume 3. CRC Press.
- Erdos, P., and Renyi, A. 1959. On random graphs. *Publ. Math. Debrecen* 6:290–297.
- Fortunato, S. 2010. Community detection in graphs. *Physics reports* 486(3-5):75–174.
- Gao, J.; Barzel, B.; and Barabási, A.-L. 2016. Universal resilience patterns in complex networks. *Nature* 530(7590):307.
- Gerstner, W.; Kistler, W. M.; Naud, R.; and Paninski, L. 2014. Neuronal dynamics: From single neurons to networks and models of cognition. Cambridge University Press.
- Han, J.; Li, Q.; et al. 2018. A mean-eld optimal control formulation of deep learning. *arXiv preprint arXiv:1807.01083*
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- Kazemi, S. M.; Goel, R.; Jain, K.; Kobayez, I.; Sethi, A.; Forsyth, P.; and Poupart, P. 2019. Relational representation learning for dynamic (knowledge) graphs: A survey. *arXiv preprint arXiv:1905.11485*
- Kingma, D. P., and Ba, J. 2015. Adam: A method for stochastic optimization. In *ICLR 2015*
- Kipf, T. N., and Welling, M. 2017. Semi-supervised classification with graph convolutional networks. *ICLR 2017*
- Kutz, J. N.; Rudy, S. H.; Alla, A.; and Brunton, S. L. 2017. Data-driven discovery of governing physical laws and their parametric dependencies in engineering, physics and biology. *2017 IEEE 7th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. IEEE.
- LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning. *Nature* 521(7553):436.
- Li, Q.; Han, Z.; and Wu, X.-M. 2018. Deeper insights into graph convolutional networks for semi-supervised learning. *Thirty-Second AAAI Conference on Artificial Intelligence*
- Lipton, Z. C.; Berkowitz, J.; and Elkan, C. 2015. A critical review of recurrent neural networks for sequence learning. *arXiv preprint arXiv:1506.00019*
- Lu, Y.; Zhong, A.; Li, Q.; and Dong, B. 2017. Beyond nite layer neural networks: Bridging deep architectures and numerical differential equations. *arXiv preprint arXiv:1710.10121*
- Luikov, A. v. 2012. Analytical heat diffusion theory. Elsevier.
- Mangan, N. M.; Brunton, S. L.; Proctor, J. L.; and Kutz, J. N. 2016. Inferring biological networks by sparse identification of nonlinear dynamics. *IEEE Transactions on Molecular, Biological and Multi-Scale Communications* 2(1):52–63.
- Narayan, A., and Roe, P. H. 2018. Learning graph dynamics using deep neural networks. *FAC-Papers Online* 5(2):433–438.
- Newman, M.; Barabasi, A.-L.; and Watts, D. J. 2010. The structure and dynamics of networks. volume 12. Princeton University Press.
- Newman, M. 2010. Networks: an introduction. Oxford U. press.
- Qin, T.; Wu, K.; and Xiu, D. 2018. Data driven governing equations approximation using deep neural networks. *arXiv preprint arXiv:1811.05537*
- Raissi, M.; Perdikaris, P.; and Karniadakis, G. E. 2018. Multistep neural networks for data-driven discovery of nonlinear dynamical systems. *arXiv preprint arXiv:1801.01236*
- Raissi, M. 2018. Deep hidden physics models: Deep learning of nonlinear partial differential equations. *The Journal of Machine Learning Research* 19(1):932–955.
- Rudy, S. H.; Brunton, S. L.; Proctor, J. L.; and Kutz, J. N. 2017. Data-driven discovery of partial differential equations. *Science Advances* 3(4):e1602614.
- Ruthotto, L., and Haber, E. 2018. Deep neural networks motivated by partial differential equations. *arXiv preprint arXiv:1804.04272*
- Seo, Y.; Defferrard, M.; Vandergheynst, P.; and Bresson, X. 2018. Structured sequence modeling with graph convolutional recurrent networks. In *International Conference on Neural Information Processing*. 362–373. Springer.
- Slotine, J.-J. E.; Li, W.; et al. 1991. Applied nonlinear control. volume 199. Prentice hall Englewood Cliffs, NJ.
- Strogatz, S. H. 2018. Nonlinear Dynamics and Chaos with Student Solutions Manual: With Applications to Physics, Biology, Chemistry, and Engineering. CRC Press.
- Thekumparampil, K. K.; Wang, C.; Oh, S.; and Li, L.-J. 2018. Attention-based graph neural network for semi-supervised learning. *arXiv preprint arXiv:1803.03735*
- Velicković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; and Bengio, Y. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903*
- Watts, D. J., and Strogatz, S. H. 1998. Collective dynamics of small-world networks. *nature* 393(6684):440.
- Wu, F.; Zhang, T.; Jr., A. H. S.; Fifty, C.; Yu, T.; and Weinberger, K. Q. 2019a. Simplifying graph convolutional networks. *ICLR*
- Wu, Z.; Pan, S.; Chen, F.; Long, G.; Zhang, C.; and Yu, P. S. 2019b. A comprehensive survey on graph neural networks. *arXiv preprint arXiv:1901.00596*
- Yang, Z.; Cohen, W. W.; and Salakhutdinov, R. 2016. Revisiting semi-supervised learning with graph embeddings. *ICML 2016* 40–48.
- Yu, B.; Yin, H.; and Zhu, Z. 2017. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. *arXiv preprint arXiv:1709.04875*
- Zang, C.; Cui, P.; Faloutsos, C.; and Zhu, W. 2018. On power law growth of social networks. *IEEE Transactions on Knowledge and Data Engineering* 30(9):1727–1740.

A Reproducibility

To ensure the reproducibility, we open-sourced our datasets and Pytorch implementation empowered by GPU and sparse matrix at:

<https://drive.google.com/open?id=19x7uas9G5w0gU8bHDoohJmDrVOt9wl8W>

B Animations of the real-world dynamics on different networks

Please view the animations of the three real-world dynamics on ve different networks learned by different models at:

<https://drive.google.com/open?id=1KBI-6Oh7BRxcQNQRPeHuKPPI6IndDa5Y>

We will find our NDCN captures the real-world dynamics on different networks very accurately while the baselines can not. The detailed experimental configurations are shown as follows:

B.1 Underlying Networks

We generate various networks by as follows, and we visualize their adjacency matrix after re-ordering their nodes by the community detection method by Newman (Newman 2010).

Grid network:

Figure 6: Adjacency matrix of grid network taking on a circulant matrix.

Random network:

```
import networkx as nx
n = 400
G = nx.erdosrenyi_graph
(n, 0.1, seed=seed)
```

Power-law network:

```
n = 400
G = nx.barabasi_albert_graph
(n, 5, seed=seed)
```

Figure 7: Adjacency matrix of random network.

Figure 8: Adjacency matrix of power-law network.

Small-world network:

```
n = 400
G = nx.newman_watts_strogatz_graph
(400, 5, 0.5, seed=seed)
```

Community network:

```
n1 = int(n/3)
n2 = int(n/3)
n3 = int(n/4)
n4 = n - n1 - n2 - n3
G = nx.random_partition_graph
([n1, n2, n3, n4], .25, .01, seed=seed)
```

B.2 Initial Values of Network Dynamics

We set the initial value (0) the same for all the experimental settings and thus different dynamics are only due to their

B.3 Network Dynamics

We adopt the following three real-world dynamics from different disciplines. Please see above animations to check out the visualization of different network dynamics. The differential equation systems are shown as follows:

The heat diffusion dynamics governed by Newton's law of cooling (Luikov 2012),

$$\frac{dx_i(t)}{dt} = k_{ij} \sum_{j=1}^n A_{ij} (x_j - x_i) \quad (6)$$

states that the rate of heat change of node i is proportional to the difference in the temperatures between it and its neighbors with heat capacity matrix A . We use $k = 1$ here.

The mutualistic interaction dynamics among species in ecology, governed by equation

$$\frac{dx_i(t)}{dt} = b_i + \sum_{j=1}^n A_{ij} \frac{x_i x_j}{d_i + e_i x_i + h_j x_j} \quad (7)$$

The mutualistic differential equation systems (Gao, Barzel, and Barabási 2016) capture the abundance $x_i(t)$ of species, consisting of incoming migration term b_i , logistic growth with population capacity k_i (Zang et al 2018) and Allee effect (Allee et al. 1949) with cold-start threshold c_i , and mutualistic interaction term with interaction network A . We use $b = 0:1$, $k = 5:0$, $c = 1:0$, $d = 5:0$, $e = 0:9$, $h = 0:1$ here.

The gene regulatory dynamics governed by Michaelis-Menten equation

$$\frac{dx_i(t)}{dt} = b_i x_i^f + \sum_{j=1}^n A_{ij} \frac{x_j^h}{x_j^h + 1} \quad (8)$$

where the first term models degradation when $f = 1$ or dimerization when $f = 2$, and the second term captures genetic activation tuned by the Hill coefficient h (Gao, Barzel, and Barabási 2016). We adopt $b = 1:0$, $f = 1:0$, $h = 2:0$ here.

B.4 Terminal Time:

We use $T = 5$ for mutualistic dynamics and gene regulatory dynamics over different networks, and $\bar{b} = 5; 0:1; 0:75; 2; 0:2$ for heat dynamics on the grid, random graph, power-law network, small-world network, and community network respectively due to their different time scale of network dynamics. Please see above animations to check out different network dynamics.

B.5 Visualizations of network dynamics

Please see above animations to check out the visualization of different network dynamics. We generate networks by aforementioned network models with $n = 400$ nodes. The nodes are re-ordered according to community detection method by Newman (Newman 2010). We visualize their adjacency matrices in Fig. 11, 12 and 13. We layout these networks in a grid and thus nodes' states are visualized as functions on the grid. Specifically, the nodes are re-ordered according to community detection method by Newman (Newman 2010) and each node has a unique label from $\{1, \dots, n\}$. We layout these nodes on a 2-dimensional $\bar{n} \times \bar{n}$ grid and each grid point $(r; c) \in \mathbb{Z}^2$ represents the i^{th} node where $i = r \bar{n} + c + 1$. Thus, nodes' states $x(t) \in \mathbb{R}^n$ when $n = 400$ can be visualized as a scalar field function $X : \mathbb{Z}^2 \rightarrow \mathbb{R}$ over the grid.

Figure 9: Adjacency matrix of small-world network.

Figure 10: Adjacency matrix of community network.

different dynamic rules and underlying networks modelled by $\dot{X} = f(X; G; W; t)$ as shown in Fig. 2,3 and 4. Please see above animations to check out different network dynamics.

```
n = 400
N = int(np.ceil(np.sqrt(n)))
x0 = torch.zeros(N, N)
x0[int(0.05 N):int(0.25 N),
    int(0.05 N):int(0.25 N)] = 25
# x0[1:5, 1:5] = 25
for N = 20 or n = 400 case
x0[int(0.45 N):int(0.75 N),
    int(0.45 N):int(0.75 N)] = 20
# x0[9:15, 9:15] = 20 for N = 20 or n = 400 case
x0[int(0.05 N):int(0.25 N),
    int(0.35 N):int(0.65 N)] = 17
# x0[1:5, 7:13] = 17 for N = 20 or n = 400 case
```

Figure 11: Heat diffusion on different networks. Each of the five vertical panels represents the dynamics on one network over physical time. For each network dynamics, we illustrate the sampled ground truth dynamics (left) and the dynamics generated by our NDCN (right) from top to bottom following the direction of time.

C Model configurations of Learning Network Dynamics in both continuous-time and regularly-sampled settings

We train our NDCN model by Adam (Kingma and Ba 2015). We choose 20 as the hidden dimension $d_h \in \mathbb{R}^n$. We train our model for a maximum of 2000 epochs using Adam (Kingma and Ba 2015) with learning rate 0.01. We summarize our regularization parameter as in Table 5 and Table 6 for Section 5 learning continuous-time network dynamics. We summarize our regularization parameter as in Table 7 for Section 6 learning regularly-sampled dynamics.

D Temporal-GNN models

We use following temporal-GNN models for structured sequence learning:

LSTM-GNN: the temporal-GNN with LSTM cell $X[t+1] = \text{LSTM}(\text{GCN}(X[t]; G))$:

$$\begin{aligned} x_t &= \text{ReLU}(W_e X[t] + b_e) \\ i_t &= \text{sigmoid}(W_{ii} x_t + b_{ii} + W_{hi} h_{t-1} + b_{hi}) \\ f_t &= \text{sigmoid}(W_{if} x_t + b_{if} + W_{hf} h_{t-1} + b_{hf}) \\ g_t &= \tanh(W_{ig} x_t + b_{ig} + W_{hg} h_{t-1} + b_{hg}) \\ o_t &= \text{sigmoid}(W_{io} x_t + b_{io} + W_{ho} h_{t-1} + b_{ho}) \\ c_t &= f_t - c_{t-1} + i_t g_t \\ h_t &= o_t \tanh(c_t) \end{aligned} \quad (9)$$

$$X[t+1] = W_d h_t + b_d$$

GRU-GNN: the temporal-GNN with GRU cell $X[t+1] =$

GRU(GCN(X[t]; G)):

$$\begin{aligned} x_t &= \text{ReLU}(W_e X[t] + b_e) \\ r_t &= \text{sigmoid}(W_{ir} x_t + b_{ir} + W_{hr} h_{t-1} + b_{hr}) \\ z_t &= \text{sigmoid}(W_{iz} x_t + b_{iz} + W_{hz} h_{t-1} + b_{hz}) \\ n_t &= \tanh(W_{in} x_t + b_{in} + r_t (W_{hn} h_{t-1} + b_{hn})) \\ h_t &= (1 - z_t) n_t + z_t h_{t-1} \\ X[t+1] &= W_d h_t + b_d \end{aligned} \quad (10)$$

RNN-GNN: the temporal-GNN with RNN cell $X[t+1] = \text{RNN}(\text{GCN}(X[t]; G))$:

$$\begin{aligned} x_t &= \text{ReLU}(W_e X[t] + b_e) \\ h_t &= \tanh(W_{ih} x_t + b_{ih} + W_{hh} h_{t-1} + b_{hh}) \end{aligned} \quad (11)$$

$$X[t+1] = W_d h_t + b_d$$

We adopt the diffusion operator $\mathcal{D} = \mathcal{D}^{-\frac{1}{2}} (I + (1 - \alpha)A) \mathcal{D}^{-\frac{1}{2}}$ where A is the adjacency matrix, \mathcal{D} is the degree matrix and $\mathcal{D} = I + (1 - \alpha)A$ keeps normalized. The differential equation system $\frac{dx}{dt} = -\mathcal{D}x$ follows the dynamics of averaging the neighborhood opinion as $\frac{dx_i(t)}{dt} = -\sum_j A_{ij} \frac{1}{(d_i + 1)} x_j(t) + x_i(t)$ for node i . When $\alpha = 0$, it averages the neighbors as normalized random walk, where \mathcal{D}^{-1} captures exponential dynamics without network effects. Here we adopt $\alpha = 0.5$, namely averages both neighbors and itself as GCN in (Kipf and Welling 2017).

E Results in absolute error.

We show corresponding loss error in Table 8, Table 9 and Table 10 with respect to the normalized loss error in Section 5

Figure 12: Biological mutualistic interaction on different networks.

Table 5: regularization parameter configurations in continuous-time extrapolation prediction

		Grid	Random	Power Law	Small World	Community
Heat	No-Encode	1e-3	1e-6	1e-3	1e-3	1e-5
	No-Graph	1e-3	1e-6	1e-3	1e-3	1e-5
Diffusion	No-Control	1e-3	1e-6	1e-3	1e-3	1e-5
	NDCN	1e-3	1e-6	1e-3	1e-3	1e-5
Mutualistic Interaction	No-Encode	1e-2	1e-4	1e-4	1e-4	1e-4
	No-Graph	1e-2	1e-4	1e-4	1e-4	1e-4
	No-Control	1e-2	1e-4	1e-4	1e-4	1e-4
	NDCN	1e-2	1e-4	1e-4	1e-4	1e-4
Gene Regulation	No-Embed	1e-4	1e-4	1e-4	1e-4	1e-4
	No-Graph	1e-4	1e-4	1e-4	1e-4	1e-4
	No-Control	1e-4	1e-4	1e-4	1e-4	1e-4
	NDCN	1e-4	1e-4	1e-4	1e-4	1e-4

learning continuous-time network dynamics and Section 6 learning regularly-sampled dynamics. The same conclusions can be made as in Table 1, Table 2 and Table 3.

F Learning Semantic labels

We summarize datasets, baselines, and experimental setups in Section 7 learning semantic labels at the terminal time.

F.1 Datasets and Baselines.

We use three standard benchmark datasets (i.e., citation network Cora, Citeseer and Pubmed), and follow the same fixed split scheme for train, validation, and test as in (Yang, Cohen, and Salakhutdinov 2016; Kipf and Welling 2017; Thekumparampil et al. 2018). We summarize the datasets in Appendix F.1 Table 11. We compare our NDCN model with graph convolution network (GCN) (Kipf and Welling 2017), attention-based graph neural network (AGNN) (Thekumparampil et al. 2018), and graph attention networks (GAT) (Velicković et al. 2017) with sophisticated attention parameters.

F.2 Experimental setup.

For the consistency of comparison with prior work, we follow the same experimental setup as (Kipf and Welling 2017; Velicković et al. 2017; Thekumparampil et al. 2018). We train our model based on the training datasets and get the accuracy of classification results from the test datasets with 1000 labels as summarized in Table 11. Following hyper-parameter settings apply to all the datasets. We set 16 evenly spaced time ticks $\{t_0, \dots, T\}$ and solve the initial value problem of integrating the differential equation systems numerically by DOPRI5 (Dormand 1996). We train our model for a maximum of 100 epochs using Adam (Kingma and Ba 2015) with learning rate 0.01 and ℓ_2 -norm regularization 0.024. We grid search the best terminal time $T \in [0.5; 1.5]$ and the $\beta \in [0; 1]$. We use 256 hidden dimension. We report the mean and standard deviation of results for 100 runs in Table 4. It's worthwhile to emphasize that in our model there is no running control parameters (i.e. linear connection layers in GNNs), no dropout (e.g., dropout rate in GCN and 0.6 in GAT), no early stop, and no concept of layer/network depth (e.g.,

Figure 13: Gene regulation dynamics on different networks.

Table 6: regularization parameter configurations in continuous-time interpolation prediction

		Grid	Random	Power Law	Small World	Community
Heat Diffusion	No-Encode	1e-3	1e-6	1e-3	1e-3	1e-5
	No-Graph	1e-3	1e-6	1e-3	1e-3	1e-5
	No-Control	1e-3	1e-6	1e-3	1e-3	1e-5
	NDCN	1e-3	1e-6	1e-3	1e-3	1e-5
Mutualistic Interaction	No-Encode	1e-2	1e-4	1e-4	1e-4	1e-4
	No-Graph	1e-2	1e-4	1e-4	1e-4	1e-4
	No-Control	1e-2	1e-4	1e-4	1e-4	1e-4
	NDCN	1e-2	1e-4	1e-4	1e-4	1e-4
Gene Regulation	No-Embed	1e-4	1e-4	1e-4	1e-4	1e-4
	No-Graph	1e-4	1e-4	1e-4	1e-4	1e-4
	No-Control	1e-4	1e-4	1e-4	1e-4	1e-4
	NDCN	1e-4	1e-4	1e-4	1e-4	1e-4

2 layers in GCN and GAT).

G Accuracy over terminal time and

By capturing the continuous-time network dynamics, our NDCN gives better classification accuracy at terminal time t^* . Indeed, when the terminal time is too small or too large, the accuracy degenerates because the features of nodes are in under-diffusion or over-diffusion states. We plot the mean accuracy of 100 runs of our NDCN model over different terminal time t and ϵ as shown in the following heatmap plots. We find for all the three datasets their accuracy curves follow rise and fall pattern around the best terminal time.

Table 7: ℓ_2 regularization parameter configurations in regularly-sampled extrapolation prediction

		Grid	Random	Power Law	Small World	Community
Heat Diffusion	LSTM-GNN	1e-3	1e-3	1e-3	1e-3	1e-3
	GRU-GNN	1e-3	1e-3	1e-3	1e-3	1e-3
	RNN-GNN	1e-3	1e-3	1e-3	1e-3	1e-3
	NDCN	1e-3	1e-6	1e-3	1e-3	1e-5
Mutualistic Interaction	LSTM-GNN	1e-3	1e-3	1e-3	1e-3	1e-3
	GRU-GNN	1e-3	1e-3	1e-3	1e-3	1e-3
	RNN-GNN	1e-3	1e-3	1e-3	1e-3	1e-3
	NDCN	1e-2	1e-3	1e-4	1e-4	1e-4
Gene Regulation	LSTM-GNN	1e-3	1e-3	1e-3	1e-3	1e-3
	GRU-GNN	1e-3	1e-3	1e-3	1e-3	1e-3
	RNN-Control	1e-3	1e-3	1e-3	1e-3	1e-3
	NDCN	1e-4	1e-4	1e-4	1e-3	1e-3

Table 8: Continuous-time Extrapolation Prediction. Our NDCN predicts different continuous-time network accurately. Each result is the ℓ_1 error with standard deviation from 20 runs for 3 dynamics on 5 networks for each method.

		Grid		Random		Power Law		Small World		Community	
Heat Diffusion	No-Encode	1:143	0:280	1:060	0:195	0:950	0:199	0:948	0:122	1:154	0:167
	No-Graph	1:166	0:066	0:223	0:049	0:260	0:020	0:410	0:023	0:926	0:116
	No-Control	2:803	0:549	1:076	0:153	0:962	0:163	1:176	0:179	1:417	0:140
	NDCN	0:158	0:047	0:163	0:060	0:187	0:020	0:097	0:016	0:183	0:039
Mutualistic Interaction	No-Encode	1:755	0:138	1:402	0:456	2:632	0:775	1:947	0:106	2:007	0:695
	No-Graph	2:174	0:089	1:038	0:434	1:301	0:551	1:936	0:085	1:323	0:204
	No-Control	5:434	0:473	1:669	0:662	9:353	3:751	4:111	0:417	2:344	0:424
	NDCN	1:038	0:181	0:584	0:277	0:653	0:230	0:521	0:124	0:502	0:210
Gene Regulation	No-Encode	2:164	0:957	6:954	5:190	3:240	0:954	1:445	0:395	8:204	3:240
	No-Graph	0:907	0:058	4:872	0:078	4:206	0:025	0:875	0:016	6:112	0:143
	No-Control	4:458	0:978	27:119	2:608	6:768	0:741	3:320	0:982	20:002	2:160
	NDCN	1:089	0:487	0:715	0:210	0:342	0:088	0:243	0:051	0:782	0:199

Figure 14: Mean classification accuracy over 100 runs of our NDCN model over terminal time and for the Cora dataset in heatmap plot.

Figure 15: Mean classification accuracy over 100 runs of our NDCN model over terminal time and for the Cora dataset in 3D surface plot.

Table 9: Continuous-time Interpolation Prediction. Our NDCN predicts different continuous-time network accurately. Each result is the ℓ_1 error with standard deviation from 20 runs for 3 dynamics on 5 networks for each method.

		Grid	Random	Power Law	Small World	Community
Heat Diffusion	No-Encode	1:222 ± 0:486	1:020 ± 0:168	0:982 ± 0:143	1:066 ± 0:280	1:336 ± 0:239
	No-Graph	1:600 ± 0:068	0:361 ± 0:022	0:694 ± 0:058	0:956 ± 0:079	0:954 ± 0:053
	No-Control	2:169 ± 0:108	1:230 ± 0:266	1:280 ± 0:216	1:544 ± 0:128	1:495 ± 0:171
	NDCN	0:121 ± 0:024	0:121 ± 0:017	0:214 ± 0:024	0:129 ± 0:017	0:165 ± 0:019
Mutualistic Interaction	No-Encode	0:620 ± 0:081	2:424 ± 0:598	1:755 ± 0:560	0:488 ± 0:077	2:777 ± 0:773
	No-Graph	0:626 ± 0:143	0:967 ± 0:269	1:180 ± 0:171	0:497 ± 0:101	1:578 ± 0:244
	No-Control	1:534 ± 0:158	2:836 ± 1:022	3:328 ± 0:314	1:212 ± 0:116	3:601 ± 0:940
	NDCN	0:164 ± 0:031	0:843 ± 0:267	0:333 ± 0:055	0:085 ± 0:014	0:852 ± 0:247
Gene Regulation	No-Encode	1:753 ± 0:555	4:278 ± 3:374	2:560 ± 0:765	1:180 ± 0:389	5:106 ± 2:420
	No-Graph	1:140 ± 0:101	3:768 ± 0:316	3:137 ± 0:264	0:672 ± 0:050	4:639 ± 0:399
	No-Control	3:010 ± 0:228	9:939 ± 1:185	3:139 ± 0:313	2:082 ± 0:293	8:659 ± 0:952
	NDCN	0:262 ± 0:046	0:455 ± 0:174	0:222 ± 0:034	0:180 ± 0:032	0:562 ± 0:130

Table 10: Regularly-sampled Extrapolation Prediction. Our NDCN predicts different structured sequences accurately. Each result is the ℓ_1 error with standard deviation from 20 runs for 3 dynamics on 5 networks for each method.

		Grid	Random	Power Law	Small World	Community
Heat Diffusion	LSTM-GNN	0:489 ± 0:081	0:824 ± 0:294	0:475 ± 0:196	0:442 ± 0:083	0:517 ± 0:162
	GRU-GNN	0:428 ± 0:085	0:349 ± 0:090	0:337 ± 0:049	0:357 ± 0:065	0:302 ± 0:031
	RNN-GNN	0:717 ± 0:227	0:957 ± 0:215	0:722 ± 0:247	0:833 ± 0:145	0:615 ± 0:000
	NDCN	0:165 ± 0:027	0:180 ± 0:063	0:208 ± 0:015	0:103 ± 0:014	0:201 ± 0:029
Mutualistic Interaction	LSTM-GNN	1:966 ± 0:126	3:749 ± 3:749	2:380 ± 0:626	2:044 ± 0:086	3:463 ± 3:095
	GRU-GNN	1:905 ± 0:157	0:162 ± 0:564	1:077 ± 0:071	1:792 ± 0:165	0:510 ± 0:549
	RNN-GNN	2:165 ± 0:004	1:303 ± 1:747	1:056 ± 0:034	2:012 ± 0:065	1:140 ± 0:887
	NDCN	1:414 ± 0:060	0:734 ± 0:168	0:990 ± 0:442	0:557 ± 0:078	0:528 ± 0:122
Gene Regulation	LSTM-GNN	1:883 ± 0:218	26:750 ± 5:634	3:733 ± 1:220	0:743 ± 0:112	16:534 ± 5:094
	GRU-GNN	1:641 ± 0:191	20:240 ± 2:549	3:381 ± 1:455	0:626 ± 0:099	14:4 ± 2:358
	RNN-GNN	1:906 ± 0:464	22:46 ± 2:276	4:036 ± 1:229	0:795 ± 0:300	14:496 ± 1:077
	NDCN	1:267 ± 0:672	0:946 ± 0:357	0:397 ± 0:133	0:312 ± 0:043	0:901 ± 0:160

Table 11: Statistics for three real-world citation network datasets. N, E, D, C represent number of nodes, edges, features, classes respectively.

Dataset	N	E	D	C	Train/Valid/Test
Cora	2,708	5,429	1,433	7	140/500/1,000
Citeseer	3,327	4,732	3,703	6	120/500/1,000
Pubmed	19,717	44,338	500	3	60/500/1,000

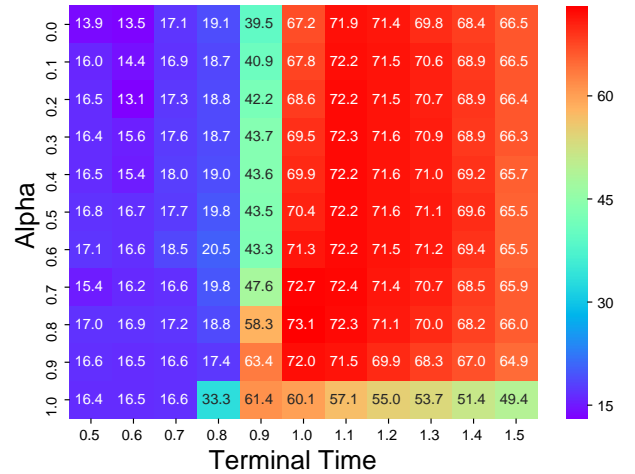


Figure 16: Mean classification accuracy of 100 runs of our NDCN model over terminal time and α for the Citeseer dataset.

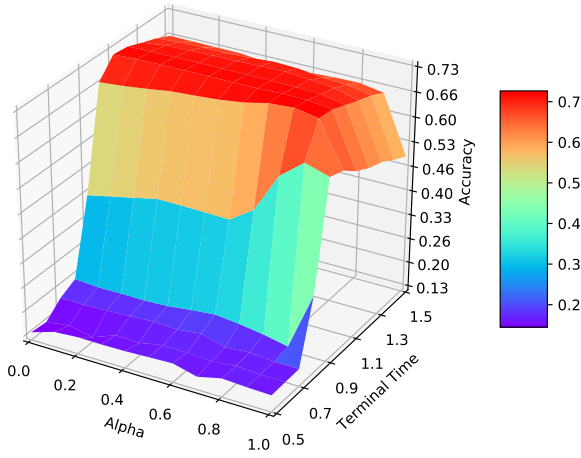


Figure 17: Mean classification accuracy of 100 runs of our NDCN model over terminal time and α for the Citeseer dataset in 3D surface plot.

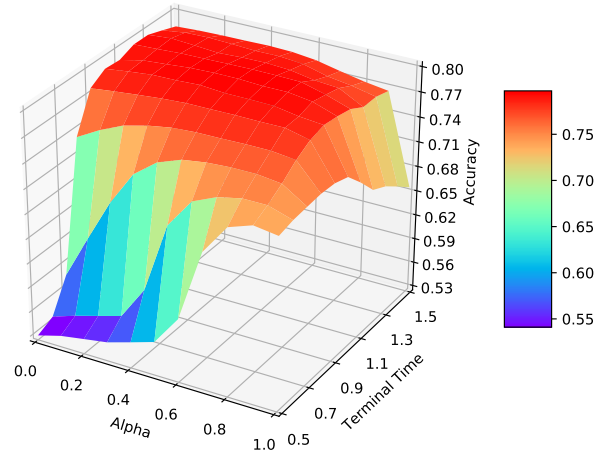


Figure 19: Mean classification accuracy of 100 runs of our NDCN model over terminal time and α for the Pubmed dataset in 3D surface plot.

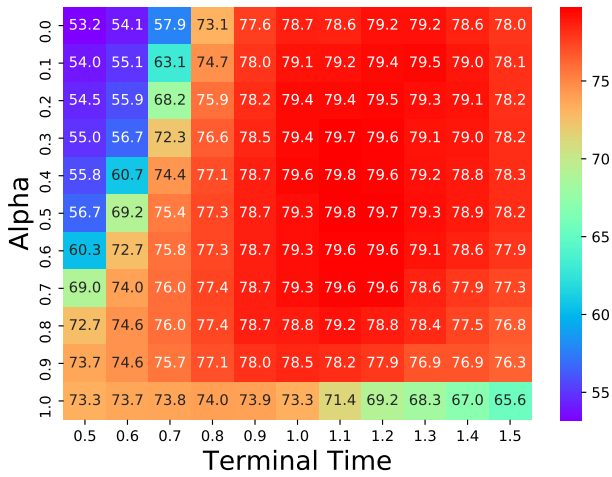


Figure 18: Mean classification accuracy of 100 runs of our NDCN model over terminal time and α for the Pubmed dataset.