# Generalizable Machine Learning in Neuroscience using Graph Neural Networks

**Paul Y Wang**[1 ✉]**, Sandalika Sapra**[2,5]**, Vivek Kurien George**[3,5]**, and Gabriel A. Silva**[3,4,5 ✉]

[1]Department of Physics, University of California San Diego
[2]Department of Electrical and Computer Engineering, University of California San Diego
[3]Department of Bioengineering, University of California San Diego
[4]Department of Neuroscience, University of California San Diego
[5]Center for Engineered Natural Intelligence, University of California San Diego

Although a number of studies have explored deep learning in neuroscience, the application of these algorithms to neural systems on a microscopic scale, i.e. parameters relevant to lower scales of organization, remains relatively novel. Motivated by advances in whole-brain imaging, we examined the performance of deep learning models on microscopic neural dynamics and resulting emergent behaviors using calcium imaging data from the nematode C. elegans. We show that neural networks perform remarkably well on both neuron-level dynamics prediction, and behavioral state classification. In addition, we compared the performance of structure agnostic neural networks and graph neural networks to investigate if graph structure can be exploited as a favourable inductive bias. To perform this experiment, we designed a graph neural network which explicitly infers relations between neurons from neural activity and leverages the inferred graph structure during computations. In our experiments, we found that graph neural networks generally outperformed structure agnostic models and excel in generalization on unseen organisms, implying a potential path to generalizable machine learning in neuroscience.

## Introduction

Constructing generalizable models in neuroscience poses a significant challenge because systems in neuroscience are typically complex in the sense that dynamical systems composed of numerous components collectively participate to produce emergent behaviors. Analyzing these systems can be difficult because they tend to be highly non-linear in how they interact, can exhibit chaotic behaviors and are high-dimensional by definition. As such, indistinguishable macroscopic states can arise from numerous unique combinations of microscopic parameters, i.e. parameters relevant to lower scales of organization. Thus, bottom-up approaches to modeling neural systems often fail since a large number of microscopic configurations can lead to the same observables (1) (2).

Because neural systems are highly degenerate and complex, their analysis is not amenable to many conventional algorithms. For example, observed correlations between individual neurons and behavioral states of an organism may not generalize to other organisms or even to repeated trials in the same individual (3) (4) (5). Hence, individual variability of neural dynamics remains poorly understood and a

fundamental obstacle to model development, as evaluation on unseen individuals often leads to subpar results. Nevertheless, neural systems exhibit universal behavior: organisms behave similarly. Motivated by the need for robust and generalizable analytical techniques, researchers recently applied tools from dynamical systems analysis to simple organisms in hopes of discovering a universal organizational principle underlying behavior. These studies, made possible by advances in whole-brain imaging, reveal that neural dynamics live on low-dimensional manifolds which map to behavioral states (6) (7). This discovery implies that although microscopic neural dynamics differ between organisms, a macroscopic/global universal framework may enable generalizable algorithms in neuroscience. Nevertheless, the need for significant hand-engineered feature extraction in these studies underscores the potential of deep learning models for scalable analysis of neural dynamics.

In this work, we examine the performance and generalizability of deep learning models applied to the neural activity of C. elegans (round worm/nematode). In particular, C. elegans is a canonical species for investigating microscopic neural dynamics because it remains the only organism whose connectome (the mapping of all 302 neurons and their synaptic connections) is completely known and well studied (8) (9) (10) (11). Furthermore, the transparent body of these worms allows for calcium imaging of whole brain neural activity which remains the only imaging technique capable of spatially resolving the dynamics of individual neurons (12). Leveraging these characteristics and insight gained from previous studies, we developed deep learning models that bridge recent advances in neuroscience and deep learning. Specifically, we first demonstrate state-of-the-art performance for classifying motor action states of C. elegans from calcium imaging data acquired in previous works. Next, we examine the generalization performance of our deep learning models on unseen worms both within the same study and in worms from a separate study published years later. We then show that graph neural networks exhibit a favourable inductive bias for analyzing both higher-order function and microscopic/neuron-level dynamics in C. elegans.

## Background

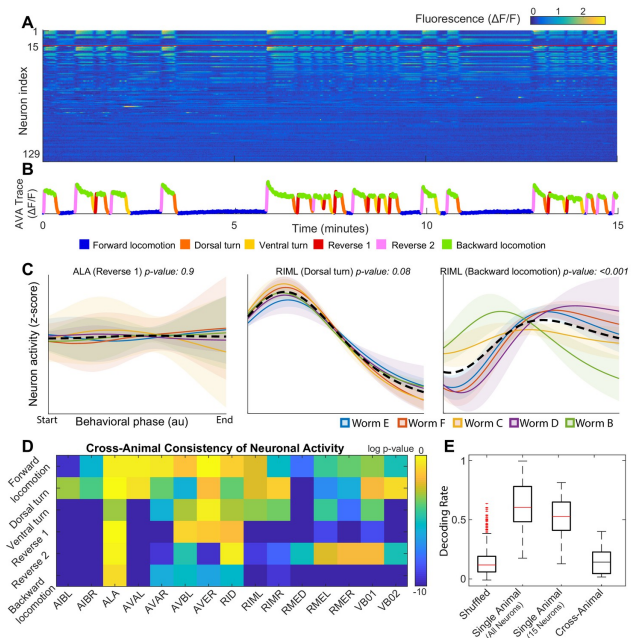In this section we discuss recent advances in neuroscience and machine learning upon which we build our model and

experiments.

**Universality/Generalizability in C. elegans models.** The motor action sequence of C. elegans is one of the only systems for which experiments on whole-brain microscopic neural activity may be performed and readily analyzed. As such, numerous efforts have focused on building models that can accurately capture the hierarchical nature of neural dynamics and resulting locomotive behaviors (13) (14). Taking advantage of this, Kato *et. al.* (7) investigated neural dynamics corresponding to a pirouette, a motor action sequence in which worms switch from forward to backward crawling, turn, and then continue forward crawling. Their analysis showed that most variations ($\sim$65%) in neural dynamics can be expressed by three principal components, and that neural dynamics in the resulting latent space trace cyclical trajectories on well-defined low dimensional manifolds corresponding to the motor action sequence (Figure S1). By identifying individual neurons, an experimental feat, these authors further determined that these topological structures in latent space were universally found among all five worms imaged in their study.
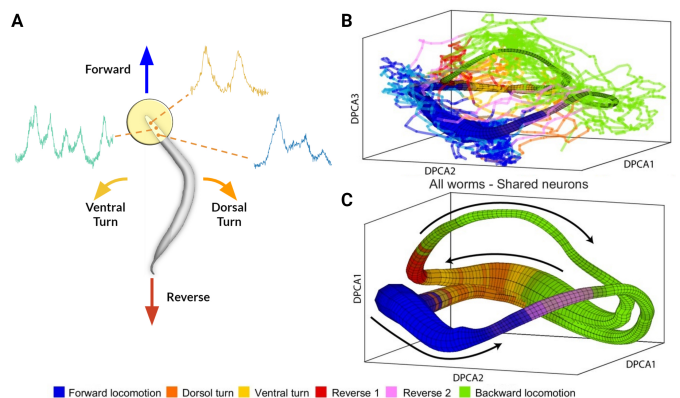
Following this initial work (7), the authors published several studies focusing on global organizational principles of C. elegans behavior (15) (16) (17). Building on two of these works, (18) found consistent differences between each individual's neural dynamics, precluding the use of established dimensional reduction techniques. For example, among 15 neurons uniquely identified among all 5 worms, only 3 neurons displayed statistically consistent behavior (Figure 1D). Examples of inconsistent behavior for unequivocally identified neurons (ALA and RIML) are shown in Figure 1C where the average of ALA's activity fails to resemble the behavior of any worm and where RIML's activity is consistent among all animals during dorsal turns, but inconsistent during reverse crawling. Resulting from these discrepancies, topological structures identified by performing PCA on each worm's neural activity were no longer observed when data from all worms was pooled together.

To address this issue, (18) introduced a new algorithm, Asymmetric Diffusion Map Modeling (ADMM), which maps the neural activity of any worm to an universal manifold (Figure 2). To achieve this, ADMM first performs time-delay embedding of neural activity into phase space. Next, a transition probability matrix is constructed by calculating distances between points in phase space using a Gaussian kernel centered on the subsequent timestep. Finally, this asymmetric diffusion map is used to construct a manifold representative of neural activity. Contrasting conventional dimensional reduction techniques, ADMM allowed quantitative modeling by mapping neural activity from the manifold, and enabled the prediction of motor action states up to 30s ahead. Despite its success, the algorithm heavily relies on hyperparameters, such as embedding parameters, which are difficult to justify and tune.

**Graph Neural Networks.** Graph Neural Networks (GNNs) are a class of neural networks that explicitly use graph structure during computations through message passing algo-



**Fig. 1.** **(A)** Calcium signals recorded in one animal for $\sim$15 minutes by (7). Each row represents a single neuron. The top 15 rows (above the red line) correspond to neurons unambiguously identified in all animals (shared neurons). **(B)** Sample trace with corresponding behavioral state colored. **(C)** Neural dynamics of two neurons for specific behavior states. Colored solid lines are the mean activity for each animal, and the black dashed line is the mean activity for all animals. Shaded colored regions show 95% confidence intervals. **(D)** Probabilities that neural dynamics from different individuals were drawn from the same distribution. **(E)** Attempt by (18) to decode onset of backwards locomotion using neural dynamics for each animal and averaged neural dynamics across other four animals. Reproduced with permission from (18).
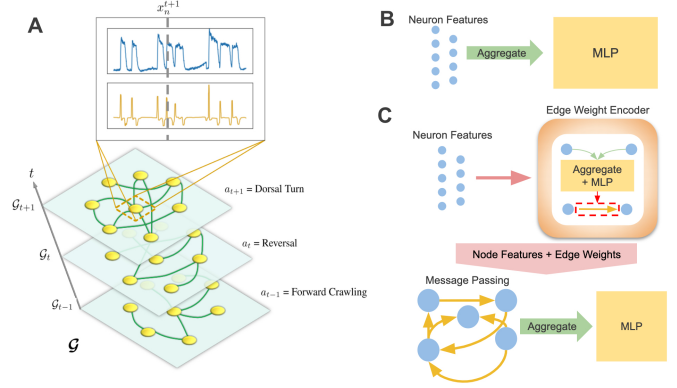


**Fig. 2.** **(A)** Rendering of calcium imaging experiment where activity of neurons in the head of the worm is recorded. Coloured arrows show main motor action behavioral states. **(B)** and **(C)** Resulting manifold from (18). **(B)** Manifold constructed from activity of four worms with coloured lines indicating neural activity of fifth worm. **(C)** Manifold constructed from neural activity of uniquely identified neurons ($n$=15) shared among all 5 worms. Black arrows correspond to cyclical transition of motor action sequence and colors correspond to motor action states. Modified with permission from (18).

rithms where features are passed along edges between nodes and then aggregated for each node ((19); (20)). These networks were inspired by the success of convolutional neural networks in the domain of two-dimensional image processing and failures when extending conventional convolutional networks to non-euclidean domains ((21)). In essence, because graphs can have arbitrary structure, the inductive bias of convolutional neural networks (equivariance to translational transformations (22)) often breaks down when applied to graphs. Addressing this issue, an early work on GNNs showed that one-hop message passing approximates spectral convolutions on graphs ((23)). Subsequent works have examined the representational power of GNNs in relation to the Weisfeiler-Lehman isomorphism test (24) and limitations of GNNs when learning graph moments ((25)). From an applied perspective, GNNs have been widely successful in a wide variety of domains including relational inference ((26); (27); (28)), node classification (23) (29)), point cloud segmentation (30), and traffic forecasting ((31) (32). In neuroscience, GNNs have been used on various tasks such as annotating cognitive state (33), and several frameworks based on graph neural networks have been proposed for analyzing fMRI data ((34); (35)).

**Relational Inference.** Relational inference remains a long-standing challenge with early works in neuroscience seeking to quantify correlations between neurons (36). Modern approaches to relational inference employ graph neural networks as their explicit reliance on graph structure forms a relational inductive bias (37) (21). In particular, our model is inspired by the Neural Relational Inference model (NRI) which uses a variational autoencoder for generating edges and a decoder for predicting trajectories of each object in a system (26). By inferring edges, the NRI model explicitly captures interactions between objects and leverages the resulting graph as an inductive bias for various machine learning tasks. This model was successfully used to predict the trajectories of coupled Kuramoto oscillators, particles connected by springs, the pick and roll play from basketball, and motion capture visualizations. Subsequently, the authors developed Amortized Causal Discovery, a framework based on the NRI model which infers causal relations from time-dependent data (27).

**Deep Learning in Neuroscience.** With the success of convolutional neural networks, researchers successfully applied deep learning to numerous domains in neuroscience ((38)) including MRI imaging (39) and connectomes (40) where algorithms can predict disorders such as autism (41). Similarly, brain-computer interfaces (BCI) are a well-studied field related to our work as they focus on decoding macroscopic variables from measurements of neural activity (42). These studies generally involve fMRI or EEG data, which characterize neural activity on a population level, to varying amounts of success (43) (44) (45) (46). Regardless, a challenge for the field is developing generalizable algorithms to individuals unseen during training (47).



**Fig. 3. (A)** Visualization of temporal graph. Inset shows $\boldsymbol{x}_n$ plotted against $t$ where the top is the calcium trace, and the bottom is its derivative. The dashed line intercepts the feature vectors at $t' = t + 1$ and denotes $x_n^{t+1}$. **(B)** and **(C)** are simplified visualizations of the MLP and GNN modules respectively.

## Model

In this section, we first present the general framework of our behavioral state classification and trajectory prediction models. Next, we detail the implementation of our neural network modules.

**Framework.** We define the set of trajectories (calcium imaging traces) for each worm as $\boldsymbol{X}_\alpha = \{\boldsymbol{x}_1, ..., \boldsymbol{x}_n\}$ where $\alpha$ denotes the label of the individual, $n$ the name of the neuron, and $\boldsymbol{x}_n$ the feature vector of the neuron. In our case, $\boldsymbol{x}_n$ corresponds to time-dependent normalized calcium traces and their derivatives for each neuron. Likewise, $x_n^t$ corresponds to the feature(s) of neuron $n$ at timestep $t$. Finally, the behavioral states of an individual are encoded as $\boldsymbol{a}_\alpha = (a^1, ..., a^t)$ where a behavioral state $a$ is assigned for each timestep $t$.

Separate models were developed for each task: behavioral state classification and trajectory prediction. In both cases, data from a worm $\alpha$ is structured as a temporal graph $\mathcal{G}_\alpha = (\mathcal{G}_\alpha^1, ..., \mathcal{G}_\alpha^t)$ (Figure 3A) where each timestep is represented by a static graph whose nodes correspond to neurons. Following the notation above, the trajectories of each neuron's calcium traces are encoded as node features $\boldsymbol{x}_n$, and the behavioral state of an individual is interpreted as a graph feature $a_\alpha^t$. For behavioral state classification, our model consists of the following (we omit $\alpha$ and $t$ in intermediary steps to simplify notation):

$$\boldsymbol{H} = f(\boldsymbol{X}_\alpha^t) \tag{1}$$

$$\boldsymbol{p} = Softmax(\boldsymbol{H}) \tag{2}$$

$$\hat{a}_\alpha^t = Max(\boldsymbol{p}) \tag{3}$$

where $f$ is an universal approximator/neural network module (described in the next section), $\boldsymbol{H}$ are hidden features, $\boldsymbol{p}$ is the probability that a system is in one of $k$ states, and $\hat{a}_\alpha^t$ is the most probable/predicted state.

For trajectory prediction, we developed a Markovian model for inferring trajectories of a consecutive timestep:

$$\boldsymbol{H} = f(\boldsymbol{X}_\alpha^t) \tag{4}$$

$$X_\alpha^{t+1} = X^t + H \qquad (5)$$

where $H$ and $f$ are the same as before. We also experimented with non-Markovian models (RNNs) for which a hidden state is included for each timestep.

The structure of our models allows us to substitute various modules for $f$. While we include results from several neural networks, we focus on two representative models: a multi-layer perceptron (MLP) agnostic to graph structure (Figure 3B) and a graph neural network (GNN) which explicitly computes on an inferred graph (Figure 3C).

**Neural Network Modules: MLP and GNN.** Our MLP module aggregates the features of a graph and feeds the aggregated features into a two-layer MLP neural network:

$$H = Aggregation(x_1^t, ..., x_n^t) \qquad (6)$$

$$H_{out}^t = g(H) \qquad (7)$$

where $g$ is a MLP. Contrasting the MLP module, our GNN relies on message passing between connected nodes and contains an encoder for edge weights $w_{ij}$:

$$H_1 = g_{enc}(X^t) \qquad (8)$$

$$H_{ij} = g(Aggregation(h_i, h_j)) \qquad (9)$$

$$p_{ij} = Softmax(h_i, h_j) \qquad (10)$$

$$w_{ij}^t = p_{ij}^2 \qquad (11)$$

where (9) encodes a hidden representation $H_{ij}$ for the edges. Applying the softmax function to $H_{ij}$ produces a two dimensional probability vector normalized to 1. We define the second dimension $p_{ij}^2$ as the weight $w_{ij}$ of an edge between nodes $i$ and $j$. The edge weights either dynamically change in each timestep's inferred graph $\mathcal{G}^t$ or remain fixed for the whole temporal graph $\mathcal{G}$ of an individual worm. If the edges are static for the temporal graph, the aggregation step in (9) also averages hidden features across all timesteps.

After edges are encoded, the GNN performs a message passing and aggregation step:

$$H_i = \sum_j^N w_{ij}^t x_j^t \qquad (12)$$

$$H_{out}^t = g(Aggregation(H_i)) \qquad (13)$$

The sum is performed over all nodes in the graph such that weighted messages are passed between connected nodes and potentially along self edges. The message passing step (12) can also be formulated in terms of an inferred weighted adjacency matrix $A^t$ and node features $X^t$:

$$H^t = A^t X^t \qquad (14)$$

Theoretically, an arbitrary number of message passing steps can be implemented; however, we did not find any improvements when using more than one step. In addition, we find that performance improves when using concatenation instead of summation during the aggregation step.

## Experiments

**Data.** Our experiments were performed with data acquired in (7) and (15). We summarize various details about the data in this section; however, we direct the reader to each respective publication for specific experimental details.

***Calcium Imaging.*** Kato *et. al.* (7) showed that neural activity corresponding to the motor action sequence lives on low dimensional manifolds. To record neuron level dynamics, they did whole-brain genetically encoded $Ca^{2+}$ imaging with single-cell-resolution and measured $\sim$100 neurons for around 18 minutes. They then normalized each calcium trace by peak fluorescence and identified neurons using spatial position and previous literature (48). Aside from imaging freely moving worms, the authors also examined robustness of topological features to sensory stimuli changes, hub neuron silencing, and immobilization. For simplicity, we limited our experiments to data collected on freely moving worms.

Nichols et. al. (15) focused on differences in neural activity of C. elegans while awake or asleep and studied two different strains of worms, n2 ($n$=11) and npr1 ($n$=10). Because experiments in both studies were performed by the same group, most experimental procedures were similar, allowing us to easily process data to match the Kato dataset. While this dataset includes imaging data of each worm during quiescence, for consistency with the Kato dataset, we only included data before sleep was induced. Furthermore, we combined results for both strains of worms as we did not notice any statistically relevant differences between them.

***Data Processing.*** We normalized the calcium trace and its derivative of each neuron to [0,1]. Normalization was performed for the entire recorded calcium trace of a worm instead of within each batch because the relative magnitudes of the traces have been found to contain graded information about the worm's behavioral state (eg. crawling speed). To create training batches, we separated each calcium trace of approximately 3000-4000 timesteps into batches of 8 timesteps where each timestep corresponds to roughly 1/3 of a second. We chose batch sizes of 8 timesteps because visualization of calcium traces showed that most local variations occur within this time frame. Moreover, 8 timesteps roughly corresponds to 3 seconds which is about the amount of time a worm needs to execute a behavioral change. Finally, the batches were shuffled before being divided into 10 folds later used for cross-validation, ensuring that each fold is representative across the whole dataset.

To compare with previous works, we performed our experiments on uniquely identified neurons between the datasets that we investigated. Identifying specific neurons is an experimental challenge, and as such, only a small fraction of neurons were unequivocally labeled. A total of 15 neurons were uniquely identified between all worms ($n$=5) measured in the Kato dataset: (AIBL, AIBR, ALA, AVAL, AVAR, AVBL, AVER, RID, RIML, RIMR, RMED, RMEL, RMER, VB01, VB02). In addition, the Nichols dataset contained data from 21 worms with 3 uniquely identified neurons

**Table 1.** Classification Accuracy of Forward and Reverse Crawling

|  | Training Set | Evaluation Set (Kato) | Evaluation Set (Nichols) |
|---|---|---|---|
| (18) | 83 | 81 | — |
| SVM | 98.8 ± .4 | 82.8 ± 7.6 | 79.0 ± 11.7 |
| MLP | 99.3 ± .6 | 93.9 ± 10.3 | 88.9 ± 11.4 |
| GNN (Connectome) | 99.5 ± .6 | 96.8 ± 4.3 | 85.5 ± 12.9 |
| GNN | **99.5 ± .5** | **97.7 ± 3.1** | **95.5 ± 6.1** |



**Fig. 4. (A)** Classification accuracy of our GNN and MLP models where black vertical lines show statistical spread. **Left**: Classification of 7 motor action states within the Kato dataset. **Right**: Classification of 4 motor action states on both the Kato and Nichols datasets. **(B)** Confusion matrix. Percent occurrence of predicted states against labeled states when evaluating on the Nichols dataset. **(C)** Mean squared error (MSE) of the GNN and various MLP models evaluated on the Nichols dataset. All models were trained using data from one worm or five worms in the Kato Dataset. **(D)** Table of MSE values for all models for 1, 8, and 16 timesteps.
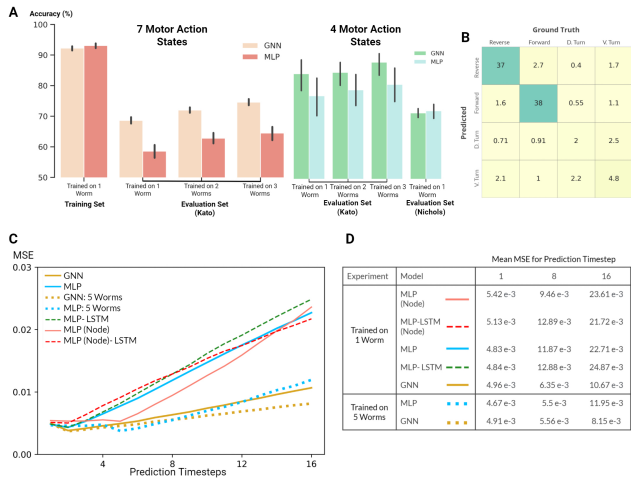
shared among all worms in both datasets: (AIBR, AVAL, VB02).

**Results.** Following (18), we used data from (7) for training/evaluating our models and data from (15) as an extended evaluation set. Because whole brain imaging is incredibly difficult, our datasets were relatively small. To address this, we experimented with data augmentation by combining data from multiple worms in the Kato dataset during model training. For all experiments, we performed 10-fold cross validation on all permutations of worms in our training set. More details, along with supplemental experiments, can be found in the Supplementary Information.

***Behavioral State Classification.*** Our first experiment compared the performance of our models to state-of-the-art results reported in (18). Specifically, this experiment involved the classification of only two motor action states, forward and reverse crawling. Along with our models described above, we also experimented with a support vector machine (SVM) and a GNN which computes with edges derived from the physical connectome (8). In particular, we incorporated the connectome into our model to investigate whether physical/structural connections between neurons can serve as a favourable inductive bias for our GNN. Our results are shown in Table 1 where Training Set denotes test set accuracy after training on the same worm and Evaluation Set denotes evaluation/generalization accuracy on worms unseen during training.

Our deep learning models clearly outperformed the SVM and state-of-the-art results, demonstrating the ability of our models to successfully classify behavioral states and generalize to other worms. Interestingly, the SVM matched the performance of our deep learning models on test set accuracy; however, its generalization performance on unseen individuals was significantly worse than our deep learning models. As such, the SVM distinctly illustrates challenges of individual variability for model development in neural systems despite the simplicity of our experiments which involve the same set of unequivocally identified neurons. Similarly, our GNN using edges derived from the connectome performed well on the test set but generalized worse than when using inferred edges. We hypothesize that the detrimental effect of using the connectome may be attributed to the model's lack of expressiveness and the distinction between inferred/functional and structural connectivity (See S.1.4.3).

Following the previous experiment, we applied our MLP and GNN models to the harder task of classifying all behavioral states labeled in the Kato dataset (Figure 4A). Within this dataset, 7 states were labeled: Forward Crawling, Forward Slowing, Reverse 1, Reverse 2, Sustained Reverse

Crawling, Dorsal Turn, and Ventral Turn. In comparison to the Kato dataset, only 4 states were labeled in the Nichols dataset: reverse crawling, forward crawling, ventral turn, and dorsal turn. For compatibility, we mapped the 7 states of the Kato dataset to 4 states of the Nichols dataset when using the Nichols dataset as an extended evaluation set.

Despite the harder task of classifying 7 states, our models achieved a classification accuracy of ∼92% on the same worm (Figure 4A: Left). Moreover, our GNN trained on three worms in the Kato dataset generalized with an accuracy of 87% (Figure 4A: Right) when classifying 4 states on the remaining unseen worms. This substantially exceeds the performance of our MLP model and (18) who report a 81% cross-animal accuracy on two states. Nevertheless, both MLP and GNN models generalized equally well (∼70%) to the 21 unseen worms of the Nichols dataset. These experiments consistently demonstrate that our GNN exceeds the performance of state-of-the-art techniques and also often exceeds the performance of our baseline MLP model.

***Neuron-level Trajectory Prediction.*** For trajectory prediction, we predicted each neuron's calcium trace and its derivative (normalized to [0,1]) for 8 timesteps during training and 16 timesteps during evaluation/validation. While training our Markovian models, scheduled sampling was performed to minimize the accumulation of error (49). In addition to our Markovian models, we also experimented with RNN implementations trained with burn-in periods of four timesteps. For evaluation, we averaged the loss per prediction timestep across all batches. Our experiments primarily focused on generalization performance of our models on the extended evaluation/Nichols dataset (Figure 4C).

Predicting neuron-level trajectory using deep learning is fairly novel since advances in whole-brain imaging are recent and limited to few organisms. Because calcium traces

are notoriously noisy and our dataset is relatively small, the performance of our model is poor; however, inspecting the MSE as a function of prediction step (Figure 4C) demonstrates that all deep learning models are able to learn transitions in the system. Moreover, increasing the number of worms included during training also improved generalization performance of our MLP and GNN models. Perhaps most surprising, our Markovian GNN outperformed all MLP models and their derived RNN variants. We attribute this result to the largely deterministic nature of neural dynamics, characterized by sparse bifurcations on the latent manifold, and the inductive bias of GNNs. As a result, given a single timestep, our GNN model was able to predict future trajectories on unseen worms for at least 16 timesteps and clearly outperformed all other models.

## Discussion

For both tasks, our GNN consistently matched or exceeded our MLP model which we accredit to its favourable inductive bias. Kato *et. al.* (7) established that projecting neural dynamics onto three principal components for each worm reveals universal topological structures; however, attempts to project neural dynamics onto shared principal components of all worms failed to display any meaningful structure. Thus, variability in each worm's neural activity, corresponding to low dimensional manifolds in latent space, is represented by different linear combinations of neurons. In other words, relevant topological structures in latent space are loosely related by linear transformations of node features. We speculate that our GNN's performance stems from its explicit structure of message passing along inferred edges which is analogous to learning linear transformations of node features (see equation (14)).

Interestingly, our model's performance was not significantly impacted by using 3 neurons ($\sim$1% of all neurons) instead of 15 ($\sim$5% of all neurons). This is not surprising because neurons strongly coupled to the motor action sequence retain most information (50), a fact consistent with (18) who found that strategically choosing 1 neuron retains $\sim$75% of the information contained in the larger set of 15 neurons.

Finally, as a critical question, we ask whether our model's performance stems from choosing a stereotyped organism that is well studied and biologically simple, or if our results imply a path towards generalizable/universal machine learning in neural systems. While the neurophysiology of C. elegans is quite complex, the motor action sequence we studied is relatively simple, especially in comparison to other organisms and cognitive functions. Moreover, organisms are adaptive and capable of learning new behavior, a fact not represented in our dataset. However, a recent astounding study (51) measured neural dynamics in monkeys trained to perform action sequences and determined that learned latent dynamics live in low-dimensional manifolds that were conserved throughout the length of the study. By aligning latent dynamics, their model accurately decoded the action of monkeys up to two years after the model was trained despite changes in biology (eg. neuron turnover, adaptation to im-

plants). Consequently, we posit that techniques similar to those used in our model may broadly apply to more complex organisms and functions.

## Conclusion

In this study, we examined the ability of neural networks to classify higher-order function and predict neuron level dynamics. In addition, inspired by global organizational principles of behavior discovered in previous studies, we demonstrated the ability of neural networks to generalize to unseen organisms. Specifically, our models exceeded the performance of previous studies in behavioral state classification of C. elegans. Furthermore, our models successfully generalized to unseen organisms, both within the same study, and in a separate experiment spaced years apart. We found that a simple MLP performs remarkably well on unseen organisms. Nevertheless, our graph neural network, which explicitly learns linear transformations of node features, matched or exceeded the performance of graph agnostic models in all experiments.

We note that our results of generalization on both higher-order functions and neuron-level dynamics (macroscopic and microscopic) suggests wide applicability of our technique to numerous machine learning tasks in neuroscience and hierarchical dynamical systems. A promising research direction is the hierarchical relationship between neuron-level and population-level dynamics. Breakthroughs in this direction may inform machine learning models working with population-level functional and imaging techniques, such as EEG or fMRI, which are readily available and widespread. In addition, in this study, we only focused on simple machine learning tasks and imaging data taken under similar experimental conditions. Further studies may involve more complex tasks such as those involving graded information in neural dynamics, changes in sensory stimuli, acquisition of learned behaviors, and higher-order functions comprised of complicated sequences of behavior. From a machine learning perspective, the development of a recurrent graph neural network with a suitable attention kernel may greatly aid model performance. Moreover, additional work is needed in examining and improving model performance on arbitrary sets of neurons as neuron identification is experimentally challenging and limited to small systems. Finally, our results show that data augmentation through the inclusion of more individuals can significantly improve generalization performance in microscopic neural systems.

## Conflict of Interest Statement

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## Author Contributions

Experiments and models were conceived by P.Y. Wang. S. Sapra assisted with the implementation of various algorithms.

The manuscript was written and revised after numerous iterations by all the authors.

## Funding

## Acknowledgments

## Data Availability Statement

The datasets (7) and (15) analyzed for this study can be found in the OSF repository here.

## Bibliography

1. Jorge Golowasch, Mark S Goldman, LF Abbott, and Eve Marder. Failure of averaging in the construction of a conductance-based neuron model. *Journal of neurophysiology*, 87(2): 1129–1131, 2002.
2. Astrid A Prinz, Dirk Bucher, and Eve Marder. Similar network activity from disparate circuit parameters. *Nature neuroscience*, 7(12):1345–1352, 2004.
3. Yves Frégnac. Big data and the industrialization of neuroscience: A safe roadmap for understanding the brain? *Science*, 358(6362):470–477, 2017.
4. Mark M Churchland, John P Cunningham, Matthew T Kaufman, Stephen I Ryu, and Krishna V Shenoy. Cortical preparatory activity: representation of movement or first cog in a dynamical machine? *Neuron*, 68(3):387–400, 2010.
5. Mark S Goldman, Jorge Golowasch, Eve Marder, and L. F. Abbott. Global structure, robustness, and modulation of neuronal models. *Journal of Neuroscience*, 21(14):5229–5238, 2001.
6. Robert Prevedel, Young-Gyu Yoon, Maximilian Hoffmann, Nikita Pak, Gordon Wetzstein, Saul Kato, Tina Schrödel, Ramesh Raskar, Manuel Zimmer, Edward S Boyden, et al. Simultaneous whole-animal 3d imaging of neuronal activity using light-field microscopy. *Nature methods*, 11(7):727–730, 2014.
7. Saul Kato, Harris S Kaplan, Tina Schrödel, Susanne Skora, Theodore H Lindsay, Eviatar Yemini, Shawn Lockery, and Manuel Zimmer. Global brain dynamics embed the motor command sequence of caenorhabditis elegans. *Cell*, 163(3):656–669, 2015.
8. John G White, Eileen Southgate, J Nichol Thomson, and Sydney Brenner. The structure of the nervous system of the nematode caenorhabditis elegans. *Philos Trans R Soc Lond B Biol Sci*, 314(1165):1–340, 1986.
9. Cornelia I Bargmann and Eve Marder. From the connectome to brain function. *Nature methods*, 10(6):483, 2013.
10. Lav R Varshney, Beth L Chen, Eric Paniagua, David H Hall, and Dmitri B Chklovskii. Structural properties of the caenorhabditis elegans neuronal network. *PLoS Comput Biol*, 7(2): e1001066, 2011.
11. Steven J Cook, Travis A Jarrell, Christopher A Brittin, Yi Wang, Adam E Bloniarz, Maksim A Yakovlev, Ken CQ Nguyen, Leo T-H Tang, Emily A Bayer, Janet S Duerr, et al. Whole-animal connectomes of both caenorhabditis elegans sexes. *Nature*, 571(7763):63–71, 2019.
12. Chentao Wen and Koutarou D Kimura. How do we know how the brain works?—analyzing whole brain activities with classic mathematical and machine learning methods. *Japanese Journal of Applied Physics*, 59(3):030501, 2020.
13. Gopal P Sarma, Chee Wai Lee, Tom Portegys, Vahid Ghayoomie, Travis Jacobs, Bradly Alicea, Matteo Cantarelli, Michael Currie, Richard C Gerkin, Shane Gingell, et al. Openworm: overview and recent advances in integrative biological simulation of caenorhabditis elegans. *Philosophical Transactions of the Royal Society B*, 373(1758):20170382, 2018.
14. Padraig Gleeson, David Lung, Radu Grosu, Ramin Hasani, and Stephen D Larson. c302: a multiscale framework for modelling the nervous system of caenorhabditis elegans. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 373(1758):20170379, 2018.
15. Annika LA Nichols, Tomáš Eichler, Richard Latham, and Manuel Zimmer. A global brain state underlies c. elegans sleep behavior. *Science*, 356(6344), 2017.
16. Harris S Kaplan, Oriana Salazar Thula, Niklas Khoss, and Manuel Zimmer. Nested neuronal dynamics orchestrate a behavioral hierarchy across timescales. *Neuron*, 105(3):562–576, 2020.
17. Susanne Skora, Fanny Mende, and Manuel Zimmer. Energy scarcity promotes a brain-wide sleep state modulated by insulin signaling in c. elegans. *Cell reports*, 22(4):953–966, 2018.
18. Connor Brennan and Alexander Proekt. A quantitative model of conserved macroscopic dynamics predicts future motor commands. *Elife*, 8:e46814, 2019.
19. F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
20. Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1263–1272, 2017.
21. Peter Battaglia, Jessica Blake Chandler Hamrick, Victor Bapst, Alvaro Sanchez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Caglar Gulcehre, Francis Song, Andy Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani, Kelsey Allen, Charles Nash, Victoria Jayne Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matt Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. *arXiv*, 2018.
22. Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999, 2016.
23. Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
24. Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.
25. Nima Dehmamy, Albert-László Barabási, and Rose Yu. Understanding the representation power of graph neural networks in learning graph topology. In *Advances in Neural Information Processing Systems*, pages 15413–15423, 2019.
26. Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural relational inference for interacting systems. In *International Conference on Machine Learning*, pages 2688–2697, 2018.
27. Sindy Löwe, David Madras, Richard Zemel, and Max Welling. Amortized causal discovery: Learning to infer causal graphs from time-series data, 2020.
28. David Raposo, Adam Santoro, David Barrett, Razvan Pascanu, Timothy Lillicrap, and Peter Battaglia. Discovering objects and their relations from entangled scene representations. *arXiv preprint arXiv:1702.05068*, 2017.
29. Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in neural information processing systems*, pages 1024–1034, 2017.
30. Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019.
31. Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. *arXiv preprint arXiv:1709.04875*, 2017.
32. Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. *arXiv preprint arXiv:1707.01926*, 2017.
33. Yu Zhang and Pierre Bellec. Functional annotation of human cognitive states using graph convolution networks. 2019.
34. Xiaoxiao Li and James Duncan. Braingnn: Interpretable brain graph neural network for fmri analysis. *bioRxiv*, 2020.
35. Byung-Hoon Kim and Jong Chul Ye. Understanding graph isomorphism network for brain mr functional connectivity analysis. *arXiv preprint arXiv:2001.03690*, 2020.
36. Clive WJ Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: journal of the Econometric Society*, pages 424–438, 1969.
37. Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, et al. Interaction networks for learning about objects, relations and physics. In *Advances in neural information processing systems*, pages 4502–4510, 2016.
38. Joshua I Glaser, Ari S Benjamin, Roozbeh Farhoodi, and Konrad P Kording. The roles of supervised machine learning in systems neuroscience. *Progress in neurobiology*, 175: 126–137, 2019.
39. Alexander Selvikvåg Lundervold and Arvid Lundervold. An overview of deep learning in medical imaging focusing on mri. *Zeitschrift für Medizinische Physik*, 29(2):102–127, 2019.
40. Colin J Brown and Ghassan Hamarneh. Machine learning on human connectome data from mri. *arXiv preprint arXiv:1611.08699*, 2016.
41. Colin J Brown, Jeremy Kawahara, and Ghassan Hamarneh. Connectome priors in deep neural networks to predict autism. In *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, pages 110–113. IEEE, 2018.
42. Gabriel A Silva. A New Frontier: The Convergence of Nanotechnology, Brain Machine Interfaces, and Artificial Intelligence. *Frontiers in Neuroscience*, 12:843, 2018. ISSN 1662-4548. doi: 10.3389/fnins.2018.00843.
43. Pouya Bashivan, Irina Rish, Mohammed Yeasin, and Noel Codella. Learning representations from eeg with deep recurrent-convolutional neural networks. *arXiv preprint arXiv:1511.06448*, 2015.
44. No-Sang Kwak, Klaus-Robert Müller, and Seong-Whan Lee. A convolutional neural network for steady state visual evoked potential classification under ambulatory environment. *PloS one*, 12(2):e0172578, 2017.
45. Arthur Mensch, Julien Mairal, Danilo Bzdok, Bertrand Thirion, and Gaël Varoquaux. Learning neural representations of human cognition across many fmri studies. In *Advances in neural information processing systems*, pages 5883–5893, 2017.
46. Joseph G Makin, David A Moses, and Edward F Chang. Machine translation of cortical activity to text with an encoder–decoder framework. Technical report, Nature Publishing Group, 2020.
47. Xiang Zhang, Lina Yao, Xianzhi Wang, Jessica Monaghan, David Mcalpine, and Yu Zhang. A survey on deep learning based brain computer interface: Recent advances and new frontiers. *arXiv preprint arXiv:1905.04149*, 2019.
48. Z. F. Altun, L. A. Herndon, C. A. Wolkow, C. Crocker, R. Lints, and D. H. Hall. Worm atlas, 2002-2020.
49. Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 1171–1179, 2015.
50. Peiran Gao and Surya Ganguli. On simplicity and complexity in the brave new world of large-scale neuroscience. *Current opinion in neurobiology*, 32:148–155, 2015.
51. Juan A Gallego, Matthew G Perich, Raeed H Chowdhury, Sara A Solla, and Lee E Miller. Long-term stability of cortical population dynamics underlying consistent behavior. *Nature neuroscience*, 23(2):260–270, 2020.
52. Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.

53. Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
54. Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493*, 2015.
55. Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. Strategies for pre-training graph neural networks. *arXiv preprint arXiv:1905.12265*, 2019.
56. Barry Horwitz. The elusive concept of brain connectivity. *Neuroimage*, 19(2):466–470, 2003.

# Supplementary Information

## Model and Experiments

**Model Selection.** The two final models included in the main text were chosen for their performance and simplicity. Nevertheless, we experimented with numerous established models which were easily substituted for $f$. For GNNs, we primarily used the Pytorch Geometric library (52). Tested modules included the GIN-0/GIN-$\epsilon$ (24), Graph Sage (29), GAT (53), and Global Attention (54). In particular, we expected the GIN to outperform the other modules because its expressiveness has been shown to aid transfer learning (55); however, because our edges are not explicitly known, we essentially applied the GIN on a fully connected graph. Under this formulation, the GIN-0 simply symmetrizes node features after a message passing step which is similar to the aggregation step of our MLP. We also found that the GIN-$\epsilon$ was prone to overfitting. Finally, we tested the GAT which is similar to our model when edges are dynamically inferred each timestep. As a result, we found that the GAT performs equally well on trajectory prediction but performs slightly worse on behavioral state classification.

**Model Implementation.** The two-layer MLP corresponding to $g$ in the main text comprised of linear layers followed by ReLu activation functions. We also applied batch norm on the output of the two layers. The Node MLP in the main text refers to individual MLPs for each node. To construct RNN variants, we added an LTSM unit before the MLP.

We performed some minor hyperparameter optimization as our combinatorial cross-validation was computationally expensive. Overall, we found our models relatively robust to different hyperparameters. For trajectory prediction, we used hidden layers with 256 dimensions. On the other hand, for behavioral state classification, we used hidden layers with 16 dimensions. Furthermore, we determined that dynamic edges evaluation worked better for trajectory prediction; however, globally evaluated edges for each worm resulted in better performance for behavioral state classification. Finally, for trajectory prediction, we chose to optimize the mean square error (MSE). For behavioral state classification, we optimized the negative log likelihood (NLL).

**Experimental Procedures.** For the extended evaluation set, we chose data from the prelethargus phase, i.e. part of the stage of larval development associated with higher frequency pharyngeal pumping prior to a cessation during which the animal enters a brief lethargus, where 4 states were labeled:

reverse, forward, dorsal turn, and ventral turn. For compatibility with the training dataset, we mapped reverse 1, reverse 2, and sustained reverse crawling to the reverse state. Similarly, we mapped forward crawling and forward slowing to forward. In addition to the 7 or 4 labeled states, there was another labeled state for unknown behavior or quiescence. This state comprised a very small portion of our data, and during training and evaluation, we ignore the result when the target is unknown.
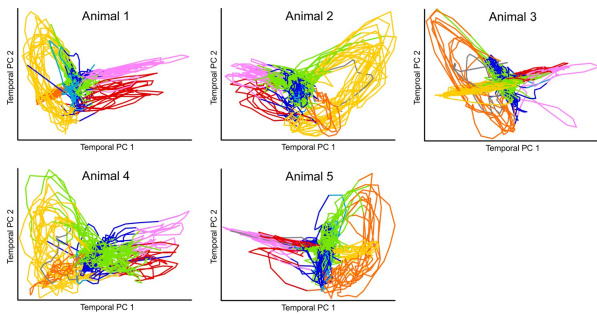
For all experiments in the main text, we performed 10-fold cross validation on all possible permutations of worms in our training set (Kato dataset). For example, on our experiments trained on two worms, the possible permutations of worms are the following: $\{(1, 2), (1, 3), (1, 4), (1, 5), (2, 3), (2, 4), (2, 5), (3, 4), (3, 5), (4, 5)\}$. Experiments labeled with "Train on 2 worms" involved models trained separately on each of these permutations. Each permutation then involved 10-fold cross validation where the test set was left out when performing hyperparamter optimization. In particular, for our experiments on behavioral state classification, we used 1 fold as the test/"leave-out" set and 1 fold for the validation set which was used for optimization and as a metric for stopping training. On the other hand, our experiments on trajectory prediction was focused primarily on generalization performance instead of test set accuracy so we used 1 fold as the validation set and evaluated on all worms in the extended validation set (Nichols dataset). As a note, we also attempted experiments where data from the extended dataset was used as a validation set. Under this condition, we found that the MLP performed significantly better; however, we were concerned that the MLP was overfitting to the validation set so we chose not to included those results.

We performed our experiments on with an Intel i9 9900k CPU and Nvidia GeForce RTX 2080Ti graphics card. Since our models are relatively simple, we were able to train the model on data from one worm in one batch. Nevertheless, the number of worms and cross-validation procedure was very computationally expensive. As such, training and evaluating each model required roughly a week or two of continuous computation. For optimization, we used the Adams optimizer with a learning rate of $10^{-3}$. We decayed the learning rate with by a factor of 0.25 if the loss did not improve after 50 epochs. We then trained for 800 epoch and saved the model with the lowest validation loss. For scheduled sampling (used during trajectory prediction), we adopted a linear decay which terminated at 300 timesteps.

**Additional Experiments.** We performed numerous experiments to verify our results and examine the performance of our model on diverse machine learning tasks. We did not perform rigorous cross validation for the following experiments.

***Experiments without AVA.*** Referees of (18) were concerned with behavioral state classification where AVA neurons were included. In particular, these neurons were used by (7) to define behavioral state through trajectory clustering in latent space. Referees commented that classifying behavioral states with neurons used to define those states was akin to circular

**Fig. S1.** Time derivatives of calcium traces projected onto each individual organism's principal components. Distinct loops correspond to manifolds in latent space where colors correspond to behavior assigned in Kato et al. Reproduced with permission from (18).

reasoning. We would like to note that (7) verified their assigned behavioral states through recorded videos, minimizing risks that assigned behavioural states differ from reality. Nevertheless, we followed (18) and performed an experiment excluding AVA neurons in which we found no noticeable difference in model performance.

***One-hot encoding of edges.*** To enforce a sparsity on the edges, we experimented with one-hot encoding by adding a scaling factor within the softmax. We found that our GNN achieved similar test accuracies as in the main text. However, our GNN failed to generalize well to unseen worms. Following our discussion in the main text, we believe that one-hot encoding was detrimental to generalization because it effectively results in a permutation matrix which simply permutes node features. This is counter to previous studies where topological structures are related by more general linear transformations.

***Comparison of inferred edges to known connectome.*** Inferring the connectivity between neurons in neural systems remains a key challenge in neuroscience. Because C. Elegans is among few organisms whose connectome mostly or completely known, we decided to compare the inferred edges of our model to the connectome of C. Elegans. Ultimately, we found no similarities between our inferred edges and the connectome.

In neuroscience, two types of connectivity are defined: structural and functional/effective. Structural connectivity refers to physical connections between neurons, whereas functional connectivity implies statistical correlations between neurons and effective connectivity validated causal connections between neurons (56). The development of methods for determining functional and, in particular, effective connectivity remains an open challenge and a highly active area of research . Nonetheless, in the context of C. elegans, each worm generally has the same structural connectivity; however, differences in neural activity implies a different functional connectivity exists for unique individuals. Since the connectome relates to the structural connectivity, we believe that our inferred edges are a poor proxy for the connectome. On a more abstract level, our graph neural network works with a subset of neurons such that a inferred edge may not correspond to a direct correlation, but may rather represent higher order correlations with unseen neurons.