
Graph Neural Stochastic Differential Equations

Richard Bergna

University of Cambridge
rsb63@cam.ac.uk

Felix L. Opolka

University of Cambridge
flo23@cam.ac.uk

Pietro Liò

University of Cambridge
pl219@cam.ac.uk

Jose Miguel Hernandez-Lobato

University of Cambridge
jmh233@cam.ac.uk

Abstract

We present a novel model *Graph Neural Stochastic Differential Equations* (Graph Neural SDEs). This technique enhances the *Graph Neural Ordinary Differential Equations* (Graph Neural ODEs) by embedding randomness into data representation using Brownian motion. This inclusion allows for the assessment of prediction uncertainty, a crucial aspect frequently missed in current models. In our framework, we spotlight the *Latent Graph Neural SDE* variant, demonstrating its effectiveness. Through empirical studies, we find that Latent Graph Neural SDEs surpass conventional models like Graph Convolutional Networks and Graph Neural ODEs, especially in confidence prediction, making them superior in handling out-of-distribution detection across both static and spatio-temporal contexts.

1 Introduction

Before the widespread use of neural networks and modern machine learning, differential equations, including Stochastic Differential Equations (SDEs) were the gold standard for modeling systems across diverse scientific disciplines [3, 4, 7, 12, 26, 29]. In the machine learning arena, neural networks’ integration with ODEs has enabled advanced continuous-time data modeling [5, 20]. Despite these progresses, the fusion of SDEs and Graph Neural Networks (GNNs) is still unexplored.

This work presents the **Graph Neural SDE**, a model that harnesses the robustness of SDEs and the versatility of GNNs to handle complex graph-structured data. Due to its stochastic nature, this model also enables precise uncertainty quantification in its predictions. Furthermore, we uncover a deep theoretical connection between our Graph Neural SDEs and the continuous representations of deep Graph Recurrent Neural Networks (refer to Appendix A). Additionally, we highlight parallels between Graph Neural ODEs and continuous deep Graph Residual Neural Networks, mirroring findings by [20] and [5].

Of particular interest, we compare the Graph Neural SDE against existing uncertainty quantification techniques for GNNs, notably the Bayesian GNN [11] and Ensemble of GNN methods [22] - two of the few established approaches for uncertainty quantification in graph data, as documented in the literature. Our experiments, shows superior performance of the Graph Neural SDE, which in many datasets surpasses both Bayesian and Ensemble approaches. Furthermore, our experiments show the Graph Neural SDE out perform Graph Neural ODEs in most spectrum of tasks, including both static and spatio-temporal datasets. The paper further details the formulation and implementation of our models, provides intuitive visualizations, and validates our approach using real-world datasets.

2 Background

2.1 Neural Ordinary Differential Equations

Neural ODEs provide an elegant approach to modeling dynamical systems using the principles of neural networks. Instead of representing data transformations as discrete layers in a traditional deep network, Neural ODEs describe them as continuous transformations parameterized by differential equations. This concept of continuous transformations within the neural network is often termed as "continuous-depth", implying that instead of having distinct layers, the network smoothly transitions and evolves data through a continuum of depths.

These transformations specify how the state of a system, denoted as $z(t)$, evolves over time. The rate of change in the state of the system is determined by a function f , which is parameterized by neural network weights. This relationship is represented as

$$\frac{dz}{dt} = f(z(t), t), \quad z(0) = z_0.$$

This differential equation implies that the state of the system at any time t is determined by accumulating the effects of the function f from the initial state z_0 up to that time. This can be articulated more explicitly as

$$z_t = z_0 + \int_0^t f(z(s), s) ds^1.$$

In the Neural ODE framework, the unknown function f , governing system dynamics, is approximated using a neural network.

2.2 Neural Stochastic Differential Equations

Stochastic Differential Equations (SDEs) have been widely employed to model real-world phenomena that exhibit randomness, such as physical systems, financial markets, population dynamics, and genetic variations [10, 31, 32]. They generalize ODEs by modeling systems that evolve continuously over time, incorporating randomness. Informally, an SDE can be seen as an ODE that integrates a certain degree of noise

$$\frac{dz}{dt} = f(z(t), t) + \epsilon(t).$$

Here, $\epsilon(t)$ represents the time-dependent noise, typically modeled using diffusion models and Brownian motion. In a more formal definition, an SDE is

$$dz(t) = \underbrace{f(t, z(t))}_{\text{drift}} dt + \underbrace{g(t, z(t))}_{\text{diffusion}} dW(t).$$

In this equation, the system state $z(t)$ at time t evolves due to two main components: the drift function f and the diffusion function g . The term $dW(t)$ denotes the infinitesimal increment of a standard Brownian motion (or Wiener process) $W(t)$, with properties like $W(0) = 0$, independent increments, and $W(t) - W(\tau)$ being normally distributed with mean 0 and variance $t - \tau$ for $0 \leq \tau < t$. The strong solution for the SDE, denoted as $z(t)$, exists and is unique under conditions where f and g are Lipschitz and $E[z(0)^2] < \infty^2$.

The **drift function** $f(t, z(t))$ represents the deterministic component of the system's evolution, describing the expected direction of change at each time point based on the current state $z(t)$.

¹In the integral expression, s is a dummy variable of integration, representing the intermediate points between 0 and t over which the function f is integrated.

²For a comprehensive and rigorous exploration of Stochastic Differential Equations, readers are referred to [17] and [30].

The **diffusion function** $g(t, z(t))$ characterizes the system’s random component, scaling the random noise introduced by the Wiener process $W(t)$, commonly known as Brownian motion.

Within the Neural SDE framework, analogous to Neural ODEs, the SDE [20] is numerically approximated. This involves evaluating the system’s response to both deterministic and random effects over time. The solution to the given SDE is represented as

$$z(t) = z(0) + \int_0^t f(s, z(s))ds + \int_0^t g(s, z(s))dW(s).$$

This equation provides an integral expression of how the state $z(t)$ evolves, subject to both deterministic (through the function f) and random influences (through the function g and the Wiener process $W(t)$).

One challenge in Neural SDEs is when the diffusion function g becomes a learnable parameter and is trained to achieve maximum likelihood, such as by directly minimizing cross-entropy or mean square error. In these situations, the diffusion function often converges to zero, turning the Neural SDE into a Neural ODE. To address this, researchers have suggested strategies like minimizing the Kullback-Leibler (KL) divergence or Wasserstein distance, forming the foundation for advanced concepts like ‘Latent SDEs’ and ‘SDE-Generative Adversarial Networks’ (SDE-GANs) [18, 20].

2.3 Graph Neural Ordinary Differential Equations

Introduced by Poli et al. [28], Graph Neural Ordinary Differential Equations (GN-ODEs) combine continuous-depth adaptability from the Neural ODEs with graph neural network structure. GN-ODEs meld the structured representation of graph data with the continuous model flexibility, providing a continuum of GNN layers. Compatible with both static and autoregressive GNN models, GDEs afford computational advantages in static contexts using the adjoint method [5] and enhance performance in dynamic situations by leveraging the geometry of the underlying dynamics.

At the heart of GN-ODEs is the representation of the dynamics between layers of GNN node feature matrices

$$\mathbf{H}(s+1) = \mathbf{H}(s) + \mathbf{F}_{\mathcal{G}}(s, \mathbf{H}(s), \theta(s)), \quad \mathbf{H}(0) = \mathbf{X}_e.$$

In this representation, \mathbf{X} denotes the initial node features of the graph, and \mathbf{X}_e is an embedding derived from various methods such as a single linear layer or another GNN layer. The function $\mathbf{F}_{\mathcal{G}}$ represents a matrix-valued nonlinear function conditioned on graph \mathcal{G} , and $\theta(s)$ is the tensor of parameters for the s -th layer. The GN-ODE model can be more succinctly expressed as:

$$\dot{\mathbf{H}}(s) = \mathbf{F}_{\mathcal{G}}(s, \mathbf{H}(s), \theta), \quad \mathbf{H}(0) = \mathbf{X}_e,$$

where s belongs to a subset \mathcal{S} of the real numbers, \mathbb{R} , typically denoted as $[t_0, t_1]$.

In the context of GN-ODEs, $\mathbf{F}_{\mathcal{G}}$ functions as a field on graph \mathcal{G} that varies with the depth or complexity of the model, which we refer to as "depth-varying". Depending on the context, the node features might be augmented to improve both computational efficiency and the model’s performance, as indicated in prior studies.

3 Graph Neural SDEs

Drawing inspiration from Graph Neural ODEs [28] and Latent SDEs [20], we introduce our methodology: Latent Graph Neural SDEs, leveraging the latent strategy for Neural SDEs.

3.1 Latent Graph Neural SDEs

Latent Graph Neural SDEs learn an latent state $z(t)$ using Graph Neural SDEs to encapsulate the data’s underlying patterns. Once determined, this state is input into a projection network f_{Ω} to generate predictions \hat{y} . The model parameterizes an Ornstein–Uhlenbeck (OU) prior process and an approximate posterior, which is another OU process.

More formally, the prior is defined by

$$d\tilde{z}(t) = f_\theta(z(t), t, \mathcal{G})dt + \sigma(\tilde{z}_t, t)dW_t,$$

where f_θ is typically set to a constant (e.g., 0 in our experiments) and σ is set to a fixed value, such as 1.0.

The approximate posterior is

$$dz(t) = f_\phi(z(t), t, \mathcal{G})dt + \sigma(z_t, t)dW_t.$$

Here, f_ϕ is parameterized by a neural network, with ϕ representing the learned weights of the network.

Both the prior and posterior drift functions, f_θ and f_ϕ respectively, dictate the dynamics of the system. $z(t)$ denotes the system state at time t , and \mathcal{G} symbolizes the graph structure. Notably, both the prior and posterior SDEs employ the same diffusion function σ but have distinct drift functions. Sharing the diffusion function ensures that the KL divergence between these processes remains finite, facilitating its estimation by sampling paths from the approximate posterior process [20]. The KL divergence between these processes is finite and can be estimated by sampling paths from the approximate posterior process.

The evidence lower bound (ELBO) is given by

$$\mathcal{L}_{ELBO}(\phi) = \mathbb{E}_{z_t} \left[\log(p(x_{t_i}|z_{t_i})) - \int_{t_0}^{t_1} \frac{1}{2} \|u(z_t, t, \phi, \theta, \mathcal{G})\|_2^2 dt \right],$$

where x_{t_i} are the observations at time t (with i in $[t_0, t_1]$), and

$$u = g(z_t, t)^{-1}[f_\phi(z_t, t, \mathcal{G}) - f_\theta(z_t, t, \mathcal{G})].$$

Given the latent state z_t , it's fed into the projection layer, f_Ω , for further prediction. The posterior predictive is then

$$p(y^*|t^*, \mathcal{G}, \mathcal{D}) = \int p(y^*|f_\Omega(z_t, t^*, \mathcal{G})) p(z|\mathcal{D}) dz \approx \frac{1}{N} \sum_{n=1}^N p(y^*|f_\Omega(z_n, t^*, \mathcal{G})).$$

As shown on the right side of this equation, the predictive distribution is approximated using Monte Carlo sampling by drawing samples z_n from the posterior $p(z|\mathcal{D})$. The variance, of the Monte Carlo mean estimation, is given by

$$\text{Var}(y) = \frac{1}{N} \sum_{n=1}^N (y_n - \bar{y})^2.$$

Furthermore, this Graph Neural SDE can be mathematically related to popular Graph Neural Networks, like Graph RNNs. For a comprehensive insight into the relationship between Graph Neural ODEs, Graph Neural SDEs, and continuous deep Residual and Recurrent Graph Neural Networks, please refer to Appendix A.

3.2 Static Dataset

We aim to predict individual voting preferences for three candidates based on their political compass and social circles, as depicted in Figure 1. The data contains inherent noise, with individuals maintaining friendships across voting preferences and their political compass not strictly dictating their voting choice, introducing randomness.

Figure 2 presents a detailed comparison across five models: GN-ODE, GCN, our Latent GN-SDE, Bayesian GCN, and Ensemble GCN — with the last one averaging predictions from five individual GCNs. The evaluation is segmented into four metrics, each depicted in its respective sub-figure.

Accuracy vs. Training Data Proportions: Training on dataset portions from 10% to 90%, the GN-SDE consistently outperformed the other models, highlighting its data efficiency.

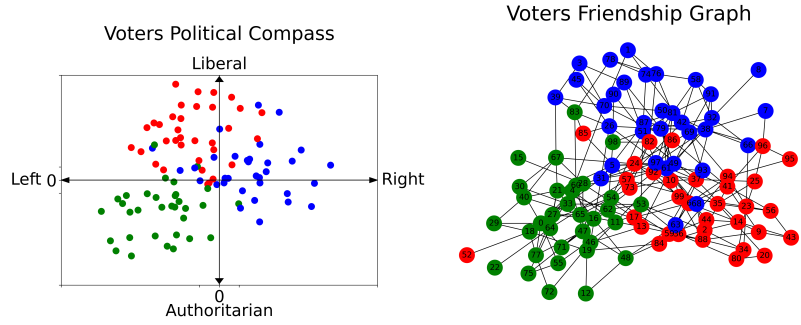


Figure 1: The left image illustrates the political compass of voters while the right image presents their social circles, with colors indicating the candidates they voted for.

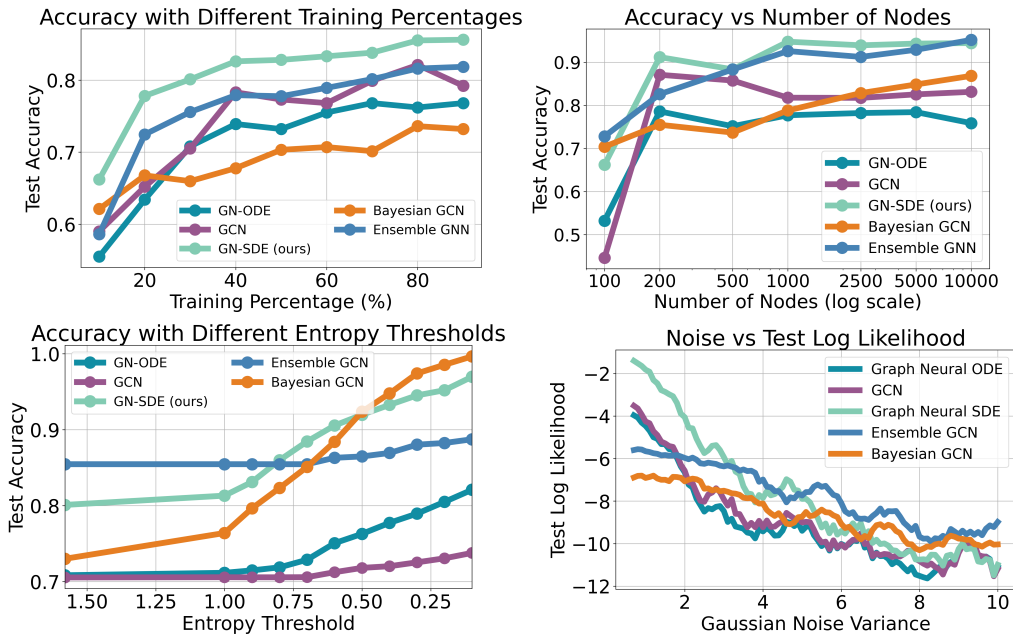


Figure 2: Comparative evaluation of Graph Neural ODE, GCN, and Graph Neural SDE models.

Accuracy vs. Number of Nodes: The GN-SDE maintained top performance across varying node counts, showcasing its adaptability to both small and large graphs. At 100 nodes, the Ensemble GCN momentarily exceeded GN-SDE, but the latter remained dominant in most scenarios.

Accuracy vs. Entropy Threshold: The entropy threshold is inversely related to the model’s confidence in its predictions; the lower the entropy, the higher the confidence required for a prediction. In this context, our model demonstrated exemplary performance in identifying out-of-distribution data and providing accurate measures of uncertainty. While our model surpassed the performance of the Ensemble GCN, it was slightly outperformed by the Bayesian GCN at very low entropy thresholds, still indicating high confidence and quality in uncertainty quantification.

Noise vs. Log-Likelihood: Evaluating resilience to added Gaussian noise, the GN-SDE model exhibited the most gradual performance decline compared to competitors like the GCN and GN-ODE. Bayesian and Ensemble models performed slightly better under high noise.

To conclude, the GN-SDE consistently surpassed GN-ODE and GCN in every test. While Bayesian GCN and the Ensemble model occasionally outperformed ours, the GN-SDE stood out in robustness, data efficiency, and uncertainty quantification.

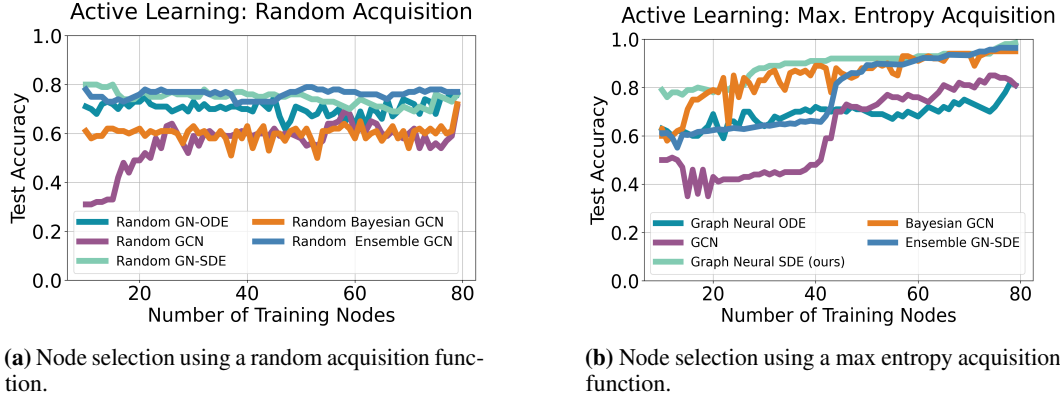


Figure 3: The figure depicts an active learning experiment on a 100-node dataset, starting with 10 nodes and incrementally adding more until reaching 80. The left and right figures use random and max entropy acquisition functions respectively for node selection.

Figure 3 depicts an active learning study on a 100-node dataset. Starting with 10 nodes, additional nodes were sequentially added: based on the highest entropy in the right figure, and randomly in the left. After each addition, the model was trained over 5 epochs, continuing until 80 nodes were included. Our model reached a notable 99% accuracy at the 78th node, closely matched by the Bayesian and Ensemble models. In comparison, GCN and Graph Neural ODE peaked around 85%. This demonstrates our model’s strength in active learning, especially vital when labeled data is limited or expensive. The random method’s ceiling at 80% accuracy highlights the benefit of uncertainty-aware strategies.

3.2.1 Real-Word Static Data sets

Dataset	Models	Entropy Thresholds											
		∞	1.6	1.0	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
CORA	GN-SDE (ours)	0.817	0.834	0.922	0.942	0.957	0.966	0.969	0.977	0.991	1.0	1.0	1.0
	GN-ODE	0.799	0.806	0.905	0.923	0.944	0.969	0.976	0.984	0.984	0.995	1.0	1.0
	GCN	0.717	0.717	0.720	0.720	0.723	0.734	0.756	0.761	0.771	0.780	0.786	0.824
	Ensemble GNN	0.777	0.802	0.935	0.949	0.954	0.958	0.962	0.972	0.983	1.0	1.0	1.0
Citeseer	GN-SDE (ours)	0.71	0.753	0.879	0.889	0.898	0.925	0.929	0.924	0.926	0.947	0.972	1.0
	GN-ODE	0.712	0.742	0.875	0.891	0.905	0.931	0.915	0.936	0.930	0.882	0.923	1.0
	GCN	0.516	0.516	0.524	0.530	0.544	0.557	0.583	0.599	0.626	0.651	0.678	0.725
	Bayesian GCN	0.61	0.619	0.729	0.746	0.769	0.791	0.815	0.821	0.854	0.863	0.867	0.88
Pubmed	GN-SDE (ours)	0.791	0.791	0.794	0.794	0.802	0.818	0.837	0.852	0.876	0.893	0.898	0.911
	GN-ODE	0.763	0.763	0.768	0.774	0.781	0.813	0.833	0.847	0.858	0.862	0.872	0.899
	GCN	0.78	0.78	0.784	0.785	0.789	0.795	0.809	0.821	0.823	0.835	0.847	0.864
	Ensemble GNN	0.786	0.786	0.796	0.814	0.837	0.868	0.879	0.908	0.907	0.916	0.929	0.952
OGB arXiv	GN-SDE (ours)	0.531	0.770	0.859	0.871	0.884	0.897	0.907	0.916	0.929	0.944	0.955	0.968
	GN-ODE	0.526	0.766	0.853	0.867	0.883	0.898	0.912	0.919	0.930	0.938	0.951	0.964
	GCN	0.470	0.809	0.885	0.900	0.909	0.917	0.931	0.939	0.950	0.966	0.983	0.976
	Ensemble GNN	0.512	0.785	0.699	0.800	1.000	1.000	1.000	-	-	-	-	-
	Bayesian GNN	0.433	0.828	0.880	0.884	0.893	0.897	1.000	0.919	0.916	0.935	0.954	1.000

Table 1: Accuracy scores for GN-SDE, GN-ODE, GCN, Ensemble GNN, and Bayesian GNN on CORA, Citeseer, Pubmed, and OGB arXiv datasets. Comparisons use varying entropy thresholds, with bold values indicating top performance. A ‘-’ for accuracy indicates the model lacked sufficiently confident data points at that threshold.

Shifting our focus to real-world datasets, Table 1 details the performance of the models GN-SDE, GN-ODE, GCN, Ensemble GCN, and Bayesian GCN on renowned datasets such as CORA [33], Pubmed [33], Citeseer [9], and OGB arXiv [13]. The models were evaluated based on different entropy thresholds. This essentially means that a model is only permitted to make predictions if its confidence level, as measured by entropy, falls below a predefined threshold. Using entropy as

a measure provides an understanding of a model’s capacity for confident predictions (uncertainty quantification) and its ability in performing out-of-distribution detection.

A review of the table shows our GN-SDE model excelling across most entropy thresholds, particularly in the CORA and Pubmed datasets. In CORA, GN-SDE achieves 100% accuracy at thresholds 0.3, 0.2, and 0.1, outperforming or matching other models. In Pubmed, GN-SDE leads across most thresholds, peaking at 91.1% accuracy for a 0.1 entropy threshold.

For Citeseer and OGB arXiv, GN-SDE’s performance is more competitive. In OGB arXiv, while GN-SDE starts strong, the Bayesian GNN reaches 100% accuracy at a 0.1 entropy threshold. Ensemble GNN also impresses, hitting 100% accuracy for thresholds up to 0.8. In Citeseer, GN-SDE often leads, with GN-ODE closely following.

In summary, GN-SDE demonstrates consistent strength across datasets, highlighting its robustness and aptitude in handling out-of-distribution data with confidence.

3.3 Spatio-Temporal

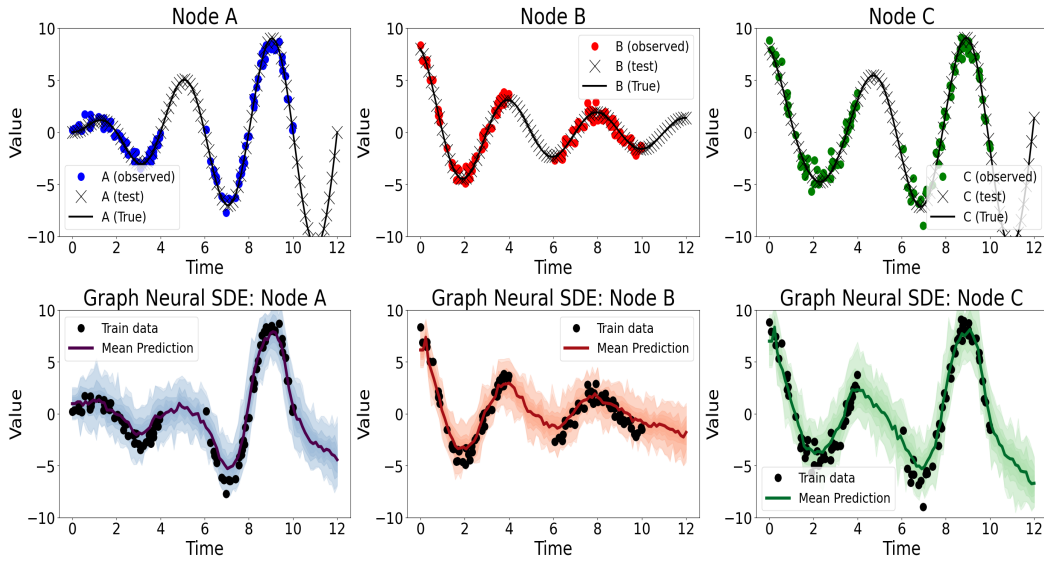


Figure 4: Comparison of the true and predicted values for the three-node regression problem. The first row shows the true values of nodes A, B, and C over time. The second row presents the predictions made by the Graph Neural SDE model for nodes A, B, and C. The shaded regions represent the model’s uncertainty quantification, demonstrating an increase in uncertainty during the interpolation and extrapolation phases.

Consider a three-node regression problem with nodes A, B, and C, where the goal is to predict their regression values at different time points. As shown in Figure 4, the observations are irregularly sampled. The graph structure used has the node C is connected to nodes A and B, but nodes A and B are not directly linked. The figure also shows the training and testing data points, illustrating both interpolation (between time 4 and 6) and extrapolation (from time 10 to 12). The true distribution for each node is as follows:

$$A(t) = t \cdot \sin\left(\frac{\pi t}{2}\right) + \epsilon_A, \quad B(t) = \frac{4}{\frac{t}{5} + 0.5} \cdot \cos\left(\frac{\pi t}{2}\right) + \epsilon_B, \quad C(t) = A(t) + B(t) + \epsilon_C.$$

where $\epsilon_A, \epsilon_B, \epsilon_C \sim \mathcal{N}(0, 0.5^2)$. It is important to note that while nodes A and B are independent, node C is a function of both A and B, with an added Gaussian noise component.

Our goal is to integrate uncertainty into predictions during testing. Both Neural ODE and GCN models are deterministic, lacking a direct uncertainty measure for regression as entropy offers

in classification. To introduce uncertainty, we use Monte Carlo Dropout during training and testing for these models. Under certain conditions, networks with dropout can emulate a Gaussian Process [8]. This allows us to estimate a model’s uncertainty by predicting both a mean and variance, using variance as the uncertainty measure. Predictions are made only if the variance meets a set threshold. The Bayesian GCN and Ensemble methods inherently quantify uncertainty, with the Ensemble calculating regression prediction variances and Bayesian GCN using Monte Carlo Sampling, aligning with our GN-SDE approach.

Metric	Models	Variance Thresholds					
		3	2.5	2	1.5	1.0	0.5
MAE	Dropout GCN	13.35	13.06	13.03	13.37	12.24	5.631
	GN-SDE (ours)	12.06	11.36	10.46	9.734	2.406	-
	Dropout GN-ODE	14.52	14.52	14.52	14.49	13.61	-
	Bayesian GCN	10.99	11.02	12.16	9.121	4.253	4.260
	Ensemble GCN	2.665	2.678	2.737	2.715	2.615	2.736
MAPE	Dropout GCN	9.340	9.390	9.462	14.08	8.941	3.060
	GN-SDE (ours)	6.825	6.321	13.08	3.380	13.55	-
	Dropout GN-ODE	8.909	8.909	8.909	10.67	9.956	-
	Bayesian GCN	7.503	7.588	9.413	7.696	7.235	2.731
	Ensemble GCN	15.44	141.20	25.15	54.49	18.31	8.612
MSE	Dropout GCN	13.35	13.06	13.03	13.37	12.24	5.631
	GN-SDE (ours)	12.06	11.36	10.46	9.734	2.41	-
	Dropout GN-ODE	14.52	14.52	14.52	14.49	13.61	-
	Bayesian GCN	10.99	11.02	12.16	9.121	4.253	4.260
	Ensemble GCN	12.73	12.72	13.27	13.28	12.55	13.47
NLL	Dropout GCN	22.55	25.07	25.29	26.74	26.66	19.81
	GN-SDE (ours)	8.969	9.032	9.530	10.30	8.44	-
	Dropout GN-ODE	40.19	40.19	40.19	42.04	43.97	-
	Bayesian GCN	13.28	13.35	12.28	12.01	6.958	9.999
	Ensemble GCN	51.38	85.56	81.54	89.95	104.68	205.54

Table 2: Performance comparison of Dropout GCN, GN-SDE, Dropout GN-ODE, and Bayesian GCN models on the three-node regression problem across different variance thresholds. Bold values indicate superior performance by the GN-SDE model.

Table 2 shows the GN-SDE model outperforming the Dropout GCN and Dropout GN-ODE models across all variance thresholds. The GN-SDE model frequently achieves the lowest Mean Squared Error (MSE) and Negative Log-Likelihood (NLL) across the higher variance thresholds. However, its performance in Mean Absolute Percentage Error (MAPE) demonstrates some variability.

For the NLL metric, calculated under the Gaussian-distributed predictions assumption

$$NLL = \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{2} \log(2\pi\sigma_i^2) + \frac{(y_i - \mu_i)^2}{2\sigma_i^2} \right),$$

where σ^2 is the predicted variance, y_i is the observations, and μ_i is the predicted mean.

As the variance threshold decreases, the Bayesian GCN becomes more competitive, even outperforming the GN-SDE in MAE at the 0.5 threshold. The Ensemble GCN’s performance fluctuates, but it remains competitive. GN-SDE’s MAE and MSE decline with decreasing variance, highlighting its capability to filter out uncertain predictions. At the 0.5 threshold, only the Dropout GCN and Bayesian GCN perform, with the latter showing remarkable MAE and NLL results.

In summary, GN-SDE is a standout in the three-node regression task, often surpassing peers across thresholds. The Bayesian GCN excels at lower variance thresholds. Dropout GCN holds its own at the strictest threshold, while Dropout GN-ODE struggles at higher variances, suggesting limitations

in uncertainty assessment. The Ensemble GCN is consistent but sometimes underperforms in the NLL metric.

To evaluate the performance of these models on a real-world dataset, MET-LA, please refer to Appendix B.2.1.

4 Related Work

Uncertainty quantification in GNNs has been relatively less explored compared to traditional neural networks. Among the limited research in this domain, we have benchmarked our work against key contributions such as the Bayesian Graph Neural Network [11] and robust ensemble methods [22]. While there has been significant research into Gaussian Processes on graphs [2, 19, 27], we did not include these in our evaluations.

The use of differential equations on graphs is a growing research area, with most work focused on Graph Neural ODEs [28] and Graph Control Differential Equations [6]. Graph Neural ODEs have been applied to dynamic graph classification [15, 34], traffic forecasting [6, 23], and protein interface prediction [36]. Extensions include second-order and higher-order Graph ODEs [25, 39].

While SDEs provide a promising avenue for better uncertainty quantification in differential equations, their application in graphs has been limited. Notably, they are employed in Graph Diffusion models primarily for denoising purposes [14, 16, 24].

In this landscape, the paper ‘BroNet’ [1] stands out for its claim to develop Graph Neural SDEs. Yet, its technique, which uses a Graph Neural Network to learn an SDE’s scalar parameter before integrating it, fails to model the graph as an SDE, setting it apart from our approach.

In essence, the field of uncertainty quantification in GNNs remains relatively unexplored. Our model presents a novel integration of SDE with GNN to address this gap. We hope that our contribution acts as a cornerstone for further developments in Graph Neural SDE-based uncertainty quantification.

5 Conclusion

This study introduces the **Graph Neural Stochastic Differential Equations** model, by using the Latent approach, designed for tasks like node, graph, and link prediction. Experimental results show that Latent Graph Neural SDEs consistently outperform models such as the Bayesian GCN, Ensemble GCN, and Graph Neural ODEs, excelling in metrics like uncertainty quantification and out-of-distribution detection. The superior accuracy of our model might stem from the intrinsic noise within its differential equations, similar to data augmentation in traditional machine learning. This noise may enhance model robustness and testing generalization, this hypothesis needs future exploration — perhaps by comparing a Graph Neural ODE trained on noise-augmented data and assessing the resulting accuracy differences. However, the advanced nature of Graph Neural SDEs requires increased computational resources, primarily due to the elevated integration costs of SDEs which make it more computationally expensive than the other models.

Moving forward, several intriguing avenues could be explored. One such avenue is the Bayesian Neural Network SDE that employs an Ornstein-Uhlenbeck prior for its weights, as described by [37]. This approach could be extended to graph SDEs, potentially leading to the development of Graph Bayesian Neural Network SDEs (Graph BNN-SDEs). Additionally, the advancements in Partial Differential Equations showcased by [35] can be creatively adapted for graph contexts using a latent method similar to our approach, paving the way for Graph Stochastic Partial Differential Equations (Graph SPDEs). Moreover, integrating higher-order stochastic differential equations with the Latent Graph Neural SDE we introduced might be advantageous and could further enhance the modeling of complex systems.

References

- [1] S. Bishnoi, S. Ranu, N. Krishnan, et al. Graph neural stochastic differential equations for learning brownian dynamics. *arXiv preprint arXiv:2306.11435*, 2023. 9

- [2] V. Borovitskiy, I. Azangulov, A. Terenin, P. Mostowsky, M. Deisenroth, and N. Durrande. Matérn gaussian processes on graphs. In *International Conference on Artificial Intelligence and Statistics*, pages 2593–2601. PMLR, 2021. 9
- [3] R. Buckdahn, B. Djehiche, and J. Li. A general stochastic maximum principle for sdes of mean-field type. *Applied Mathematics & Optimization*, 64(2):197–216, 2011. 1
- [4] L. Cardelli. From processes to odes by chemistry. In *Fifth Ifip International Conference On Theoretical Computer Science–Tcs 2008*, pages 261–281. Springer, 2008. 1
- [5] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018. 1, 3
- [6] J. Choi, H. Choi, J. Hwang, and N. Park. Graph neural controlled differential equations for traffic forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 6367–6374, 2022. 9
- [7] M. Cvijovic, J. Almquist, J. Hagmar, S. Hohmann, H.-M. Kaltenbach, E. Klipp, M. Krantz, P. Mendes, S. Nelander, J. Nielsen, et al. Bridging the gaps in systems biology. *Molecular Genetics and Genomics*, 289:727–734, 2014. 1
- [8] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016. 8
- [9] C. L. Giles, K. D. Bollacker, and S. Lawrence. Citeseer: An automatic citation indexing system. *Proceedings of the third ACM conference on Digital libraries*, pages 89–98, 1998. 6
- [10] V. Gontis and A. Kononovicius. Consentaneous agent-based and stochastic model of the financial markets. *PLoS One*, 9(7):e102201, 2014. 2
- [11] A. Hasanzadeh, E. Hajiramezani, S. Boluki, M. Zhou, N. Duffield, K. Narayanan, and X. Qian. Bayesian graph neural networks with adaptive connection sampling. In *International conference on machine learning*, pages 4094–4104. PMLR, 2020. 1, 9
- [12] S. Hoops, R. Hontecillas, V. Abedi, A. Leber, C. Philipson, A. Carbo, and J. Bassaganya-Riera. Ordinary differential equations (odes) based modeling. In *Computational Immunology*, pages 63–78. Elsevier, 2016. 1
- [13] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems*, 33:22118–22133, 2020. 6
- [14] H. Huang, L. Sun, B. Du, Y. Fu, and W. Lv. Graphgdp: Generative diffusion processes for permutation invariant graph generation. In *2022 IEEE International Conference on Data Mining (ICDM)*, pages 201–210. IEEE, 2022. 9
- [15] M. Jin, Y. Zheng, Y.-F. Li, S. Chen, B. Yang, and S. Pan. Multivariate time series forecasting with dynamic graph neural odes. *IEEE Transactions on Knowledge and Data Engineering*, 2022. 9
- [16] J. Jo, S. Lee, and S. J. Hwang. Score-based generative modeling of graphs via the system of stochastic differential equations. In *International Conference on Machine Learning*, pages 10362–10383. PMLR, 2022. 9
- [17] R. Khasminskii. *Stochastic Stability of Differential Equations*. Springer, 2012. 2
- [18] P. Kidger, J. Foster, X. C. Li, and T. Lyons. Efficient and accurate gradients for neural sdes. *Advances in Neural Information Processing Systems*, 34:18747–18761, 2021. 3
- [19] N. Lawrence and M. Jordan. Semi-supervised learning via gaussian processes. *Advances in neural information processing systems*, 17, 2004. 9
- [20] X. Li, T.-K. L. Wong, R. T. Chen, and D. Duvenaud. Scalable gradients for stochastic differential equations. In *International Conference on Artificial Intelligence and Statistics*, pages 3870–3882. PMLR, 2020. 1, 3, 4
- [21] Y. Li, R. Yu, C. Shahabi, and Y. Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *International Conference on Learning Representations (ICLR '18)*, 2018. 14
- [22] Q. Lin, S. Yu, K. Sun, W. Zhao, O. Alfarraj, A. Tolba, and F. Xia. Robust graph neural networks via ensemble learning. *Mathematics*, 10(8):1300, 2022. 1, 9

- [23] Z. Liu, P. Shojaee, and C. K. Reddy. Graph-based multi-ode neural networks for spatio-temporal traffic forecasting. *arXiv preprint arXiv:2305.18687*, 2023. 9
- [24] T. Luo, Z. Mo, and S. J. Pan. Fast graph generative model via spectral diffusion. *arXiv preprint arXiv:2211.08892*, 2022. 9
- [25] X. Luo, J. Yuan, Z. Huang, H. Jiang, Y. Qin, W. Ju, M. Zhang, and Y. Sun. Hope: High-order graph ode for modeling interacting dynamics. 2023. 9
- [26] V. Mandelzweig and F. Tabakin. Quasilinearization approach to nonlinear problems in physics with application to nonlinear odes. *Computer Physics Communications*, 141(2):268–281, 2001. 1
- [27] F. Pérez-Cruz, S. Van Vaerenbergh, J. J. Murillo-Fuentes, M. Lázaro-Gredilla, and I. Santamaria. Gaussian processes for nonlinear signal processing: An overview of recent advances. *IEEE Signal Processing Magazine*, 30(4):40–50, 2013. 9
- [28] M. Poli, S. Massaroli, J. Park, A. Yamashita, H. Asama, and J. Park. Graph neural ordinary differential equations. *arXiv preprint arXiv:1911.07532*, 2019. 3, 9
- [29] M. Quach, N. Brunel, and F. d’Alché Buc. Estimating parameters and hidden variables in non-linear state-space models based on odes for biological networks inference. *Bioinformatics*, 23(23):3209–3216, 2007. 1
- [30] D. Revuz and M. Yor. *Continuous martingales and Brownian motion*, volume 293. Springer Science & Business Media, 2013. 2
- [31] L. M. Ricciardi. *Biomathematics and related computational problems*. Springer Science & Business Media, 2012. 2
- [32] S. J. Schreiber, M. Benaïm, and K. A. Atchadé. Persistence in fluctuating environments. *Journal of Mathematical Biology*, 62:655–683, 2011. 2
- [33] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad. Collective classification in network data. In *AI magazine*, volume 29, page 93, 2008. 6
- [34] G. Shi, D. Zhang, M. Jin, and S. Pan. Towards complex dynamic physics system simulation with graph neural odes. *arXiv preprint arXiv:2305.12334*, 2023. 9
- [35] Y. Sun, L. Zhang, and H. Schaeffer. Neupde: Neural network based ordinary and partial differential equations for modeling time-dependent data. In *Mathematical and Scientific Machine Learning*, pages 352–372. PMLR, 2020. 9
- [36] X. W. Tan. Fuzzy3dvecnet: Novel robust deep learning approach to protein-ligand interaction prediction. 9
- [37] W. Xu, R. T. Chen, X. Li, and D. Duvenaud. Infinitely deep bayesian neural networks with stochastic differential equations. In *International Conference on Artificial Intelligence and Statistics*, pages 721–738. PMLR, 2022. 9
- [38] B. Yu, H. Yin, and Z. Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. *arXiv preprint arXiv:1709.04875*, 2017. 14
- [39] Y. Zhang, S. Gao, J. Pei, and H. Huang. Improving social network embedding via new second-order continuous graph neural networks. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 2515–2523, 2022. 9

A What Are Graph Neural ODE and Graph Neural SDEs anyways?

Here, we demonstrate the equivalences or representations of Graph Neural ODEs and Graph Neural SDEs in relation to popular architectures.

A.1 Graph Neural ODEs as Continuously-deep Graph Residual Neural Networks

In this section, we present a straightforward extension of this idea to graph structures. Specifically, we illustrate that a Graph Neural ODE can be conceptualized as a continuous-depth version of a Graph Residual Network.

Considering the architecture of a residual graph network:

$$y_{j+1} = y_j + f_{\mathcal{G}}(j, y_j, \theta) \quad (1)$$

where $f_{\Theta}(j, y_j, \mathcal{G})$ represents the j -th residual block, with the parameters from all blocks being collectively represented by Θ .

In contrast, let's look at the Graph Neural ODE (abbreviated as GN-ODE):

$$\frac{dy}{dt}(t) = f_{\mathcal{G}}(j, y_j, \theta)$$

Discretizing this GN-ODE using the explicit Euler method at uniformly spaced time intervals t_j with a gap of Δt gives:

$$\frac{y(t_{j+1}) - y(t_j)}{\Delta t} \approx \frac{dy}{dt}(t_j) = f_{\mathcal{G}}(j, y_j, \theta)$$

Simplifying, we get:

$$y(t_{j+1}) = y(t_j) + \Delta t \cdot f_{\mathcal{G}}(j, y_j, \theta)$$

By integrating the factor of Δt into $f_{\mathcal{G}}$, this equation naturally aligns with the formulation in Equation 1. Such a perspective underscores that neural ODEs can be seen as the continuous-time counterparts of residual networks.

This viewpoint casts the graph neural ODE as a continuously-deep Graph Residual Network. Here, the sequence of minor (residual) updates to its hidden states become both infinitesimally small and infinitely frequent. The end output is the cumulative effect of these continuous updates, mirroring the solution to the ODE from its initial state.

A.2 Graph Neural SDEs as Continuously-Deep Recurrent Graph Neural Networks

A well-established analogy exists between numerically discretized neural stochastic differential equations and the structures found in deep learning literature - particularly Recurrent Neural Networks (RNNs). In the case of Neural SDEs, the RNN's input can be interpreted as random noise, or Brownian motion, while its output corresponds to a generated sample.

Consider an autonomous one-dimensional Itô Stochastic Differential Equation represented as:

$$dy(t) = f(y(t))dt + \sigma(y(t))dw(t)$$

where $y(t)$, $f(y(t))$, $\sigma(y(t))$, and $w(t)$ belong to the set of real numbers, \mathbb{R} . The numerical Euler-Maruyama discretization of this SDE can be expressed as:

$$\frac{y(t_{j+1}) - y(t_j)}{\Delta t} \approx f(y(t_j)) + \frac{\sigma(y(t_j))\Delta w_j}{\Delta t}$$

Which simplifies to:

$$y_{j+1} = y_j + f(y_j)\Delta t + \sigma(y_j)\Delta w_j$$

Here, Δt represents a fixed time step and Δw_j is normally distributed with mean zero and variance Δt . This numerical discretization is reminiscent of an RNN with a specific form. Therefore, we can consider the Neural SDE as an consciously-deep RNN where the depth is defined by the number of discretization steps of the SDE solver.

B Extended Results

B.1 Spatial-Temporal Images

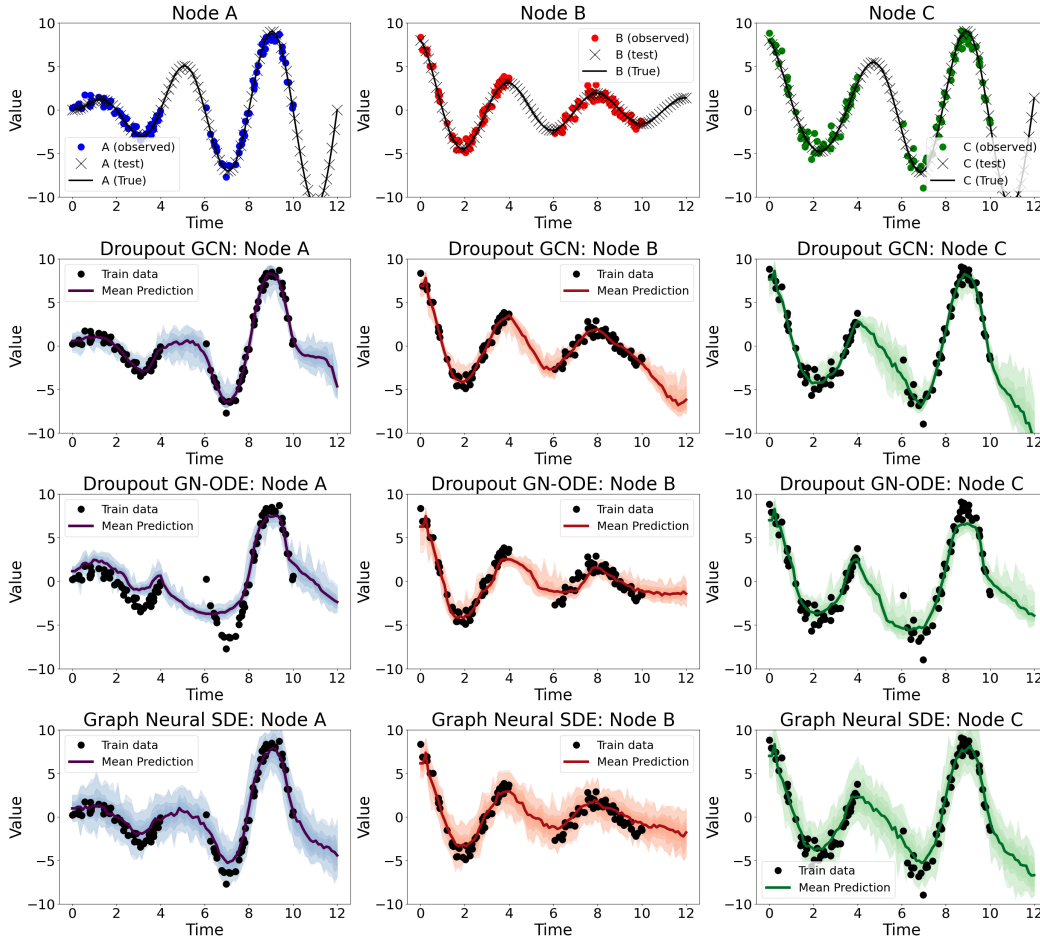


Figure 5: The figure displays 12 images organized. Each column corresponds to one of the 3 nodes, while each row represents a different model or dataset. The top row illustrates the training and testing datasets for each node, in the context of node regression. The aim is to predict the regression value for each node. The second row presents results from the GCN, the third row showcases those from the Graph Neural ODE, and the bottom row depicts our model, the Graph Neural SDE.

B.2 Meter-LA Experiment Set Up

In our experiments, we utilized Graph Attention Networks (GAT) to embed the input data, which consisted of the past six recordings. These inputs were embedded into 64-dimensional tensors. Subsequently, these embeddings were passed to our Latent Graph Stochastic Differential Equation (SDE) model. The Graph Latent SDE model employs a GCN for the drift function, while the diffusion function is held constant at a value of 1. The hidden state of the model is of size 64, which is then passed to a GCN projection layer for prediction.

The Graph Neural ODE model follows a similar structure but replaces the SDE with an

ODE and omits the noise component. The GCN model also utilizes the same embedding and projection layers but bypasses the differential equations entirely.

This setup allows for a fair comparison between the models, as they share the same embedding and projection layers, differing only in the differential equation component.

B.2.1 Real World Datasets

Transitioning to real-world spatiotemporal datasets, we focus on METR-LA for traffic prediction, a typical time-series problem. The goal is to forecast future traffic metrics, such as speed or flow, over the next H steps based on previous M steps’ traffic observations [21]. This can be mathematically described as:

$$\hat{v}_{t+1}, \dots, \hat{v}_{t+H} = \arg \max_{v_{t+1}, \dots, v_{t+H}} \log P(v_{t+1}, \dots, v_{t+H} | v_{t-M+1}, \dots, v_t),$$

where $v_t \in \mathbb{R}^n$ is an observation vector at time t for n road segments [38]. The METR-LA dataset, collected from 207 highway loop detectors in LA County, contains four months of data from March to June 2012. Data was recorded every 5 minutes, and our experiments aim to predict an hour’s worth of traffic speed (12 readings) based on the last hour’s data. Please refer to appendix B.2 for experimental setup details.

Metric	Models	Variance Thresholds				
		100	3	1	0.5	0.25
MAE	Dropout GAT-GCN	14.30	14.01	9.38	5.25	5.22
	GN-SDE (ours)	14.13	12.42	10.14	5.40	4.82
	Dropout GN-ODE	15.21	15.53	12.36	9.37	5.58
	Bayesian GCN	15.0	13.9	9.4	5.3	5.2
	Ensemble GCN	9.835	6.665	2.615	2.736	-
MAPE	Dropout GAT-GCN	10.21	10.59	20.42	21.01	20.87
	GN-SDE (ours)	10.63	9.66	18.74	18.79	13.25
	Dropout GN-ODE	10.27	13.92	21.86	17.14	13.70
	Bayesian GCN	10.5	10.6	20.5	20.9	20.8
	Ensemble GCN	18.15	15.44	18.31	8.61	-
RMSE	Dropout GAT-GCN	18.58	16.39	9.58	5.26	5.22
	GN-SDE (ours)	15.38	13.17	10.12	5.22	4.78
	Dropout GN-ODE	19.73	17.32	16.68	12.87	7.33
	Bayesian GCN	19.0	16.5	9.6	5.3	5.2
	Ensemble GCN	19.74	12.73	7.55	6.47	-
RNLL	Dropout GAT-GCN	7.03	13.28	15.12	11.56	14.52
	GN-SDE (ours)	6.56	11.57	13.54	13.92	13.88
	Dropout GN-ODE	32.33	75.10	31.03	92.42	213.09
	Bayesian GCN	7.2	13.4	15.3	12.0	14.6
	Ensemble GCN	39.63	51.38	104.68	205.54	-

Table 3: Comparative performance of various models on the METR-LA Dataset across different variance thresholds. Bold values indicate superior performance by the model.

In Table 3, the GN-SDE model consistently surpasses other contenders across all variance thresholds, notably securing the lowest MAE, MAPE, RMSE, and RNLL scores. This emphasizes GN-SDE’s stellar accuracy and unmatched uncertainty estimation abilities. Even with stricter variance thresholds of 0.5 and 0.25, GN-SDE remains at the forefront, signifying its effective prediction alignment with observed data and adeptness at uncertainty quantification.

Specifically, the Dropout GN-ODE struggles with the RNLL metric, hinting at a subpar fit and diminished uncertainty evaluation capabilities. This shortcoming becomes glaringly apparent at the 0.25 variance threshold, where its RNLL significantly eclipses that of its peers. Bayesian GCN

deserves mention for its resilient performance, often ranking a close second to the GN-SDE, underlining its robust modeling and uncertainty estimation prowess. Meanwhile, Ensemble GCN displays intermittent strengths in the MAE at lower variances but demonstrates inconsistency, particularly with fluctuating RNLL scores at the 100 and 3 variance thresholds.

In summation, while GN-SDE dominates in its robustness across varying thresholds, Bayesian GCN also emerges as a formidable model. Dropout GAT-GCN and Dropout GN-ODE especially grapple with the RNLL metric, and Ensemble GCN exhibits variable results.