

---

# Interpretable Regime Trajectories via Generative Graph State-Space Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Forecasting the behavior of real-world spatiotemporal systems often requires not  
2 only accurate predictions but also interpretable regime trajectories, i.e. discrete  
3 states that describe how dynamics change over time. However, existing approaches  
4 often entangle space and time, obscuring regime structure or trading interpretability  
5 for scale. We introduce ReGraSS, a unified framework that learns discrete,  
6 interpretable latent regimes from spatiotemporal data, represented as dynamic  
7 graphs, combining variational training with strictly time-ordered state-space inference.  
8 Predictions are produced by a mixture-of-experts modulated by the inferred  
9 regime probabilities, enforcing regime-specific specialization and supporting interpretability.  
10 Trained with self-supervised one-step prediction, the model learns in label-scarce  
11 settings and provides calibrated uncertainty by estimating a distribution  
12 over discrete regimes. ReGraSS matches or surpasses state-of-the-art spatiotemporal  
13 baselines in one-step forecasting. It shows the smallest error spike at regime  
14 changes and the fastest recovery thereafter, indicating regime-level interpretability  
15 and reliable trajectory tracking without compromising accuracy. We believe  
16 our interpretable, uncertainty-aware framework for regime-aware forecasting on  
17 dynamic graphs has direct application in healthcare, finance, and epidemiology.

## 18 1 Introduction

19 Recent advances in modern sensing and data acquisition reveal how real-world systems evolve across  
20 space and time [1, 2, 3]. In these settings, accurate forecasting of the system’s trajectory is necessary  
21 but often not sufficient: interpretable regime trajectories, i.e. discrete states governing the system’s  
22 evolution, may also be required. They reveal to be critical in high-stakes domains such as healthcare  
23 (e.g., disease progression staging [4, 5, 6]), finance (e.g., market regimes [7, 8]) or epidemiology (e.g.,  
24 transmission phases [9, 10]), where decisions rely on understanding when and why a regime shifts,  
25 not only on accurate forecasting of future events. Yet, current learning systems lack regime-aware,  
26 time-ordered explanations alongside forecasts, leaving a critical gap for models that jointly learn  
27 spatial structure, temporal evolution, and discrete interpretable regimes.

28 Graph Neural Networks (GNNs) [11, 12] are powerful tools to model spatial relationships through  
29 relational inductive biases, providing a unified framework for domains with hierarchical structure and  
30 rich spatial interactions. Extending them to spatiotemporal settings is challenging because both the  
31 graph topology and the node signals may evolve over time. Dynamic variants such as EvolveGCN  
32 [13] and ROLAND [14] address part of this challenge, yet they still interleave spatial aggregation  
33 with temporal updates via recursive message passing. As topology evolves, this coupling obscures  
34 what changed from when it changed, hindering interpretable identification of discrete regimes.

35 On the other hand, State Space Models (SSMs) provide a well-established framework for modeling  
36 temporal dynamics via latent state representations and structured transitions. Recent work on

37 learnable SSMs (e.g., S4 [15], Mamba [16]), achieve strong performance on sequence modeling tasks,  
38 overcoming several limitations of classical SSMs such as adaptation to regime shifts and multi-scale  
39 dynamics. However, they typically require large amounts of data and tend to sacrifice latent-state  
40 interpretability, limiting their applicability where understanding the underlying dynamics is essential.

41 Recent efforts integrate SSMs with graphs by decoupling spatial and temporal reasoning, applying  
42 state-space updates at the node level and mixing via GNN layers (e.g., GrassNet [17], Graph Mamba  
43 [18]). Dynamic variants further interleave Mamba-based sequence modules and spatiotemporal graph  
44 blocks to handle evolving topologies (STG-Mamba [19], DG-Mamba [20]). However, these models  
45 typically use SSMs as feature extractors rather than for interpretable regime tracking through state  
46 representation.

47 In this work, we propose **ReGraSS**, an unified framework that models discrete, interpretable regimes  
48 as latent states on dynamic graphs. A dynamic GNN encodes evolving structure and features, while  
49 temporal regime dynamics are decoupled and captured via a variational distribution over latent  
50 states. Predictions are state-conditioned via a mixture-of-experts weighted by regime probabilities, so  
51 the inferred state actively governs emissions, yielding trajectory-level explanations and calibrated  
52 uncertainty without degrading forecast accuracy. During training, a categorical VAE [21] with  
53 a learnable transition prior supports uncertainty-aware regime discovery; at inference, we switch  
54 to a strictly time-ordered state-space rollout that conditions only on past and present, enabling  
55 transparent trajectory analysis without future leakage. To function in label-scarce settings common in  
56 high-stakes applications, we adopt an autoregressive one-step forecasting objective that forces the  
57 model to internalize graph-coupled dynamics by predicting next-step node features and produces  
58 regime trajectories consistent with predictive performance. On controlled synthetic tests with induced  
59 non-stationarity, the framework captures regime transitions, supports uncertainty-aware trajectories,  
60 and matches or surpasses strong spatio-temporal and graph-SSM baselines, demonstrating that  
61 interpretable regime tracking can be achieved without a trade-off in accuracy.

## 62 2 Proposed Approach

### 63 2.1 Problem Statement

64 We consider a discrete-time sequence of graph snapshots  $\mathcal{G} = (G_t)_{t=0}^T$ , where each  $G_t = (V_t, E_t, X_t)$   
65 consists of a vertex set  $V_t$ , an edge set  $E_t$ , and node features  $X_t \in \mathbb{R}^{N_t \times D}$  with  $N_t = |V_t|$  and  
66 feature dimension  $D$ . We assume that the topology and vertex set may vary over time (vertices may  
67 appear or disappear).

68 Our main hypothesis is that a finite set of discrete regimes  $\mathcal{R} = \{r_1, \dots, r_K\}$  modulates the dynamics,  
69 with  $R_t \in \mathcal{R}$  being the active regime at time  $t$ . Discrete regimes align with how practitioners typically  
70 characterize system progression (physiological stages, market states) even when these categories  
71 coarsen underlying continuous dynamics.

72 Focusing on node-feature dynamics, we assume that the next time point features  $X_{t+1}$  are gen-  
73 erated from some probability distribution  $\mathbb{P}(X_{t+1} \mid X_{0:t}, V_{0:t}, E_{0:t}, R_{0:t})$  conditioned on past  
74 history. Thus our objective is to learn a model  $\hat{f}$  that (i) approximates the probability distri-  
75 bution  $\mathbb{P}(X_{t+1} \mid X_{0:t}, V_{0:t}, E_{0:t}, R_{0:t})$  in an autoregressive manner, (ii) while inferring the ac-  
76 tive regime without being provided regime annotations. Formally, the model can be defined as  
77  $\hat{f}(X_{0:t}, V_{0:t}, E_{0:t}) = (\hat{X}_{t+1}, \hat{R}_t)$ . Extensions to topology prediction are straightforward.

### 78 2.2 Architecture

79 We introduce **ReGraSS** (Regime-aware Graph State Space model), an autoregressive generative  
80 framework for modeling spatio-temporal dynamics on graphs through an interpretable discrete latent  
81 state space, where the extracted states act as proxies for the underlying regimes governing the system  
82 evolution. ReGraSS follows a structured encoder-decoder design. The encoder approximates the  
83 posterior over discrete latent states and the decoder generates node features at future time steps,  
84 conditioned on both the latent state and observed inputs. The model operates differently during  
85 training and inference; we describe the training behavior here and defer the dual representation and  
86 inference details to Section 2.4. The architecture is illustrated in figure 2, and we describe its main  
87 building blocks below.

Figure 1: Visualization of the model’s architecture and dual formulation.

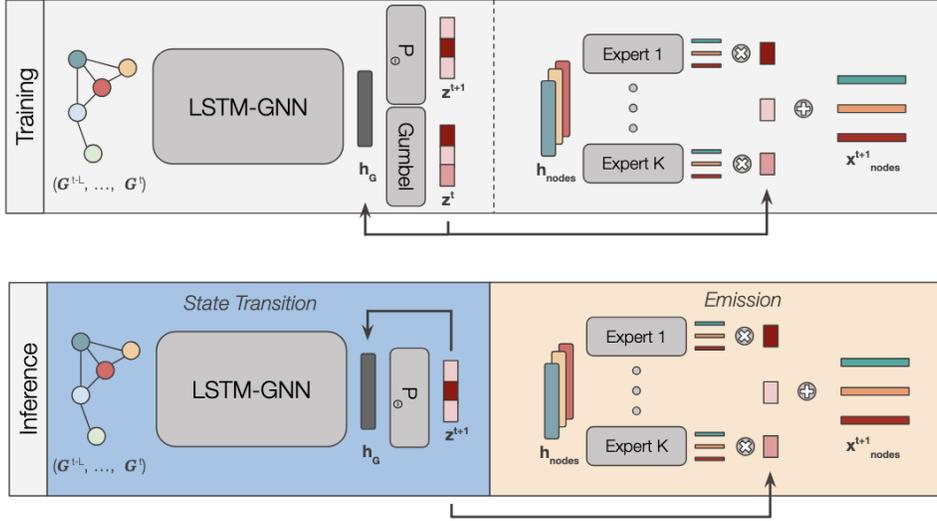


Figure 2: Visualization of the model architecture and dual-view formulation. **Top panel** (training phase): At each time step, graph snapshots are processed sequentially through the LSTM-GNN module to produce graph-level temporal embeddings. These embeddings are used to infer the current latent state  $z_t$  via the Gumbel-Softmax module, and to predict the next state  $z_{t+1}$  via the learnable prior module  $P_\theta$ . Given the inferred states, the temporal node-level embeddings  $h_{\text{nodes}}$  are passed through a MoE module, where expert outputs are modulated by  $z_t$ , to generate the predicted node features at time  $t + 1$ . **Bottom panel** (inference structure): The framework can be decomposed into two components: state transition and emission. This mirrors the classical state-space model (SSM) formulation, while extending it to a non-linear and graph-based setting.

88 **Encoder.** The encoder’s first stage is a temporal GNN that aggregates information from past snapshots  
 89 up to the current step  $t$ . We instantiate it with ROLAND ([14]), which maintains hierarchical  
 90 representations via GRU updates ([22]) and naturally supports evolving graph topology. We map the  
 91 pooled temporal graph embedding  $h_G^t \in \mathbb{R}^H$  to  $K$  unnormalized logits with a linear layer

$$\ell_t = Wh_G^t + b, \quad W \in \mathbb{R}^{K \times H}, b \in \mathbb{R}^K, K = |\mathcal{R}|.$$

92 To obtain a posterior over the  $K$  regimes, we use the Gumbel–Softmax reparameterization [21]:

$$q_\phi(z_t | h_G^t) = \text{softmax}\left(\frac{\ell_t + g_t}{\tau}\right), \quad g_t \sim \text{Gumbel}(0, 1)^K, \tau > 0.$$

93 Sampling  $z_t \sim q_\phi(\cdot | h_G^t)$  yields a differentiable, discrete latent vector that encodes the current  
 94 regime  $R_t$  (approaching one-hot as  $\tau \rightarrow 0$ ). Our probabilistic approach quantifies uncertainty in  
 95 the current regime and offers a distributional view that bridges continuous dynamics and discrete  
 96 regimes. Yet, the framework remains compatible with continuous latents if the underlying system is  
 97 better described by continuous variables.

98 **Learnable prior for causal transitions.** We define a learnable prior  $p_\theta(z_t | h_G^{t-1}, z_{t-1})$ ,  
 99 parametrized by a 2-layer MLP, that receives the temporal embedding  $h_G^{t-1}$  and the previous latent  
 100 state to predict  $z_t$ . During training, we align this prior with the variational posterior  $q_\phi$  (see Sec-  
 101 tion 2.3), yielding SSM-like transitions and enabling the dual representation in Section 2.4. Compared  
 102 with a fixed Markov prior, this data-driven conditional prior better captures non-stationary regime  
 103 dynamics on evolving graphs.

104 **Decoder (mixture of experts).** The decoder predicts the next-step node features  $X_{t+1}$  conditioned  
 105 on the current features  $X_t$  and the latent state  $z_t$ . We implement it as a mixture-of-experts (MoE  
 106 [23]):  $K$  experts  $\{f_k\}_{k=1}^K$ , each parametrized by an independent MLP, produce candidate outputs  
 107 that are combined using the posterior mixing coefficients  $\pi_t = q_\phi(z_t | h_G^t)$ . Equivalently,  $\hat{X}_{t+1} =$   
 108  $\sum_{k=1}^M \pi_{t,k} f_k(X_t)$ . The choice of the number of experts is domain specific but is typically selected to

109 match the number of regimes in the data ( $K = |\mathcal{R}|$ ), thus encouraging a one-to-one correspondence  
 110 between the variational-induced states  $z_t$  and the regimes  $r_k$ . This implementation induces state-  
 111 conditioned output generation, mirroring classical SSM two-stage behavior, i.e., state transition then  
 112 output emission [15]. It also encourages specialization: as the Gumbel–Softmax vector approaches  
 113 one-hot, each expert learns the dynamics associated with a specific regime  $r_k \in \mathcal{R}$ . While our  
 114 implementation uses MLPs, experts can be replaced with other modules such as GNNs when domain  
 115 requires it, e.g. when regime transitions influence the diffusion process in the graph, which is better  
 116 captured by GNN experts rather than MLPs.

### 117 2.3 Training Procedure

118 Training uses a variational objective derived from the categorical VAE ELBO ([21],[24]) with a  
 119 forecasting likelihood. Formally, with  $q_z^t = q_\phi(z_t | h_G^t)$  and  $p_z^t = p_\theta(z_t | h_G^{t-1}, z_{t-1})$ , the graph-level  
 120 loss over a sequence  $t = 0, \dots, T - 1$  is

$$\mathcal{L} = \sum_{t=0}^{T-1} \left[ \underbrace{\ell_{\text{forecast}}(\hat{X}_{t+1}, X_{t+1})}_{\text{one-step prediction}} + \beta \underbrace{(\text{KL}(q_z^t) \parallel \text{sg}[p_z^t])}_{\text{encoder-prior alignment}} + \gamma \underbrace{\text{CE}(\text{sg}[q_z^t], p_z^t)}_{\text{teacher-forced prior fitting}} \right], \quad (1)$$

121 where  $\ell_{\text{forecast}}$  is a regression loss between the decoder prediction  $\hat{X}_{t+1}$  (the MoE output) and the  
 122 observed features  $X_{t+1}$ , and  $\text{sg}[\cdot]$  denotes the stop-gradient operator, i.e. that the gradients are not  
 123 backpropagated further in the computation tree. The KL term updates the encoder so that the posterior  
 124  $q_z^t$  agrees with the prior  $p_z^t$ , while the cross-entropy (CE) term trains the transition module to match  
 125 the encoder’s next-time posterior  $q_z^{t+1}$ . This asymmetric pairing stabilises learning: the encoder does  
 126 not chase a moving prior, and the prior learns from the encoder without backpropagating through its  
 127 inputs. Detailed regularization and parameters schedules are provided in the Appendix 5.1.

128 By regressing  $X_{t+1}$  from information available at time  $t$ , the objective forces the model to internalize  
 129 the system’s transition mechanisms, remaining effective when regime annotations are missing,  
 130 unreliable, or available only at endpoints. This, in turn, enables reconstruction of regime trajectories  
 131 and stratification of sequences  $(G_t)_{t=0}^T$  by regime and temporal evolution.

### 132 2.4 Dual Representation

133 Our framework couples variational training with state-space inference to bridge two limitations  
 134 encountered in the literature. By learning spatial representations with a dynamic GNN and evolving  
 135 them through discrete regimes, it disentangles space–time updates that obscure regime structure  
 136 in temporal GNNs. At the same time, the inference-time state-space rollout restores interpretable  
 137 state representation often lost in deep SSMs, while preserving forecasting accuracy through state-  
 138 conditioned emissions. This dual formulation is robust to scarce or unreliable labels and preserves  
 139 strict temporal causality. During **training**, a variational next-step regression objective learns a  
 140 posterior over regime trajectories, enabling trajectory-level explanations even without ground-truth  
 141 regime annotations. At **inference**, we replace the posterior, that benefits from future information  
 142 via backpropagation, with the learned transition prior and roll forward using only past observations,  
 143 eliminating future leakage. The probabilistic treatment yields calibrated uncertainty for the current  
 144 regime and a distributional view of transitions, bridging continuous dynamics and discrete regimes  
 145 without sacrificing predictive performance.

## 146 3 Experiments and Results

### 147 3.1 Dataset

148 We generate 150 spatiotemporal graph sequences with  $T = 10$  snapshots (TP1–TP10). At TP1,  
 149 we sample  $C \sim \text{Unif}\{3, 4, 5\}$  Gaussian clusters in  $\mathbb{R}^d$  ( $d = 8$ ), with centers  $\mu_c \sim \text{Unif}([-10, 10]^d)$ ,  
 150 isotropic covariance  $0.6^2 I_d$ , and sizes  $M_c \sim \text{Unif}\{5, \dots, 100\}$ ; node features are the sampled  
 151 coordinates. We build an undirected  $k$ -NN graph at TP1 and keep edges fixed thereafter (translation-  
 152 invariant under our dynamics). Each sequence follows a discrete regime  $r_t \in \{r_1, r_2, r_3\}$  that induces  
 153 a constant drift  $v(r_t) \in \{-2 \mathbf{1}_d, +2 \mathbf{1}_d, \mathbf{0}_d\}$ , with i.i.d. Gaussian noise  $\varepsilon_i^{(t)} \sim \mathcal{N}(0, 0.6^2 I_d)$  at each  
 154 step. The regime is resampled once between TP4 and TP5 to test models ability to remain robust to a  
 155 mid-sequence nonstationarity (TP4→TP5). Full details are in Appendix 5.2.

156 **3.2 Baseline Methods**

157 We compare against baselines spanning complementary assumptions: (i) no space/no time, (ii)  
 158 time-only naïve dynamics, (iii) spatio-temporal without latent regimes, and (iv) spatio-temporal with  
 159 deep state-space modules, to ensure gains are not attributable to unstructured aggregation or trivial  
 160 autocorrelation. **MLP** (no space, no time) concatenates all node features into a single embedding,  
 161 testing whether simple global aggregation suffices. **Persistence** (time only) is a parameter-free  
 162 baseline that predicts  $X_{t+1} = X_t$  to assess whether autocorrelation alone explains performance.  
 163 **LSTM-GNN** (spatio-temporal) uses the GNN-LSTM encoder (adapted from ROLAND [14]) as a  
 164 standalone predictor, isolating the contribution of discrete regimes and mixture-of-experts decoding  
 165 in our method. **STG-MAMBA ([19])** (spatio-temporal) integrates dynamic graph filtering with a  
 166 Mamba block for multi-scale temporal modeling, providing a benchmark against state-space/dynamic-  
 167 graph hybrids. These baselines rule out unstructured aggregation, trivial autocorrelation, and generic  
 168 spatio-temporal encodings. Additional implementations details appear in the Appendix 5.3.

169 **3.3 One-Step Prediction under Changing Regimes**

170 First, we evaluate each model in a one-step regression setup to test whether it captures system  
 171 evolution and adapts to regime changes. Given the observed history up to time  $t$ ,  $(X_{0:t}, V_{0:t}, E_{0:t})$ ,  
 172 each model predicts the next features  $X_{t+1}$ ; we apply this procedure iteratively across time points on  
 173 the dataset in Section 3.1. We pay particular attention to the induced shift between TP4 and TP5 as  
 174 a stress test for non-stationarity. Performance is quantified by mean squared error (MSE) between  
 175  $\hat{X}_{t+1}$  and  $X_{t+1}$ , and we additionally report the mean absolute feature value at each time point to  
 176 contextualize error magnitude (Table 1).

177 Across the synthetic dataset, ReGraSS attains the lowest mean error over the sequence (2.79 MSE  
 178 vs. 3.05 for LSTM-GNN; Table 1) and leads both before the induced shift (pre-TP5 average  
 179 1.33) and after it (post-TP5 average 3.96). The Persistence baseline ( $x_{t+1} = x_t$ ) performs worst  
 180 throughout (6.77 MSE), confirming that temporal autocorrelation alone does not explain performance.  
 181 A structure-free MLP is competitive early but breaks at the shift (TP5), indicating that unstructured  
 182 aggregation cannot adapt to non-stationarity. The LSTM-GNN (ROLAND-based [14]) is a strong  
 183 spatio-temporal encoder without latent regimes; it matches or narrowly beats our method at isolated  
 184 time points (TP4 and TP9), yet falls behind on average and recovers more slowly after the shift.  
 185 STG-Mamba ([19]) underperforms on this setting, especially near the regime change, suggesting  
 186 limited robustness to non-stationary dynamics.

187 Two observations highlight the intended advantages of discrete regimes with state-conditioned  
 188 emissions. First, the performance drop at the regime change is the smallest for our method (TP5-  
 189 TP4 jump 4.72 vs. 4.99 for LSTM-GNN, 5.36 for MLP, 5.02 for STG-Mamba), indicating better  
 190 alignment to the new dynamics. Second, our method shows the fastest one-step recovery (TP5-TP6  
 191 drop  $-3.71$  vs.  $-3.40$  for LSTM-GNN) and post regime changes performances, consistent with rapid  
 192 state reassignment and expert specialization once the system switches regimes. A residual limitation  
 193 is a mild degradation within long single-regime segments (e.g., TP4 and TP9), which we attribute  
 194 to occasional regime misassignment due to insufficient penalty on remaining in an incorrect state  
 195 (see Figure 3). This suggests a simple mitigation with a calibrated self-transition regularizer without  
 196 altering the overall architecture.

Table 1: Validation performance on synthetic data at different time-points (TP1-TP9) of the sequence  $(G_t)_{t=1}^9$ . For each method and time step we compute the MSE between  $\hat{X}_{t+1}$  and  $X_{t+1}$ .

Model	TP1	TP2	TP3	TP4	TP5	TP6	TP7	TP8	TP9
Persistence	6.57	6.55	6.69	6.66	6.90	6.92	6.92	6.97	6.74
MLP	2.26	1.98	2.17	1.87	7.23	3.79	3.46	4.12	3.92
LSTM-GNN	1.21	1.46	2.15	<b>1.65</b>	6.64	3.24	3.55	3.86	<b>3.72</b>
STG-Mamba	8.59	6.97	3.78	3.48	8.50	5.11	4.89	4.32	4.65
Our Method	<b>1.01</b>	<b>1.14</b>	<b>1.40</b>	1.75	<b>6.47</b>	<b>2.76</b>	<b>3.11</b>	<b>3.51</b>	3.96
Features Mean	3.52	3.98	4.76	5.77	5.96	6.27	6.80	7.56	7.26

197 **3.4 Regime Trajectory Analysis**

198 To evaluate the framework’s ability to recover latent dynamics without supervision, we visualize the  
 199 learned latent trajectories. We evaluate (i) unsupervised recovery of state trajectories and (ii) the  
 200 speed of convergence to an identifiable latent-state distribution. Regime estimation was performed by  
 201 sampling the learned posterior distribution 100 times per graph and time step, followed by majority  
 202 voting to assign discrete state labels to regimes. This setup enables us to track how the inferred state  
 203 distribution evolves over time, and how it aligns with ground-truth regimes.

204 In Figure 3, we visualize inferred state trajectories in our unsupervised setting. The model initially  
 205 fails to recover the true state distribution. This is expected, as dynamic patterns must be inferred from  
 206 sequential observations alone as no regime-predictive features are present. However, after a few time  
 207 steps, the model converges to the true underlying regime distribution, demonstrating its capacity to  
 208 infer system dynamics without regime-level supervision.

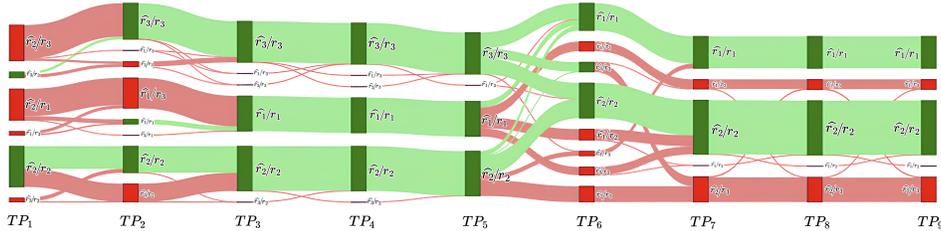


Figure 3: Unsupervised recovery of regime trajectories. The flow diagram shows one trained run across time points. Each block represents the distribution of samples (patients) by predicted/ground-truth regime pair, labeled  $\hat{r}_i/r_j$ . Green indicates correct assignments ( $\hat{r}_i = r_j$ ); red indicates mismatches. Line widths encode the number of samples flowing between pairs over time. After a short transient the mass concentrates on correct pairing flows, showing that the model recovers the latent regimes and tracks their dynamics without supervision.

209 **4 Conclusion**

210 We proposed a generative framework that integrates dynamic graph neural networks with discrete  
 211 state space modeling to capture interpretable spatio-temporal dynamics. By separating spatial  
 212 reasoning (via GNNs) from temporal inference (via a discrete latent state and learnable transition  
 213 prior), our approach addresses key limitations of prior DGNN and deep SSM models, namely limited  
 214 interpretability, entangled updates, and challenges in modeling evolving graph structures through  
 215 discrete regimes changes. The variational training procedure enables uncertainty-aware learning  
 216 of state transitions and current regime estimation, while the inference-time state-space formulation  
 217 supports forecasting without future leakage and trajectory analysis. Across synthetic experiments,  
 218 ReGraSS achieves competitive predictive accuracy while exposing latent change in regimes aligned  
 219 with system dynamics. Results show that our framework can recover temporal regimes from minimal  
 220 supervision, highlighting its utility in settings with sparse labels.

221 Our study has limitations that point to potential next directions. First, the current objective emphasizes  
 222 feature dynamics and may underweight structural change in the graph; incorporating topology-aware  
 223 terms could better capture evolving connectivity, though care is needed to avoid prohibitive costs on  
 224 large graphs. Second, the mixture-of-experts decoder scales with the number of discrete regimes,  
 225 which can hinder efficiency in fine-grained settings; lighter parameter-sharing schemes may retain  
 226 state-conditioned emissions with lower overhead. Third, a purely discrete latent space can be rigid  
 227 when regimes overlap or evolve smoothly. Beyond methodology, our evaluation on controlled  
 228 synthetic sequences should be complemented by real-world deployments that test its capabilities to  
 229 maintain interpretable regimes tracking in setups with noisy samples and complex spatio-temporal  
 230 dynamics.

## References

- 231
- 232 [1] Longqi Liu, Ao Chen, Yuxiang Li, Jan Mulder, Holger Heyn, and Xun Xu. Spatiotemporal  
233 omics for biology and medicine. *Cell*, 187(17):4488–4519, 2024.
- 234 [2] Ali Hamdi, Khaled Shaban, Abdelkarim Erradi, Amr Mohamed, Shakila Khan Rumi, and  
235 Flora D Salim. Spatiotemporal data mining: a survey on challenges and open problems.  
236 *Artificial Intelligence Review*, 55(2):1441–1488, 2022.
- 237 [3] Mariana Belgiu and Alfred Stein. Spatiotemporal image fusion in remote sensing. *Remote  
238 sensing*, 11(7):818, 2019.
- 239 [4] Ronnachai Jaroensri, Ellery Wulczyn, Narayan Hegde, Trissia Brown, Isabelle Flament-Auvigne,  
240 Fraser Tan, Yuannan Cai, Kunal Nagpal, Emad A Rakha, David J Dabbs, et al. Deep learning  
241 models for histologic grading of breast cancer and association with disease prognosis. *NPJ  
242 breast cancer*, 8(1):113, 2022.
- 243 [5] Ling Dai, Liang Wu, Huating Li, Chun Cai, Qiang Wu, Hongyu Kong, Ruhan Liu, Xiangning  
244 Wang, Xuhong Hou, Yuexing Liu, et al. A deep learning system for detecting diabetic retinopathy  
245 across the disease spectrum. *Nature communications*, 12(1):3242, 2021.
- 246 [6] Shuai Niu, Jing Ma, Qing Yin, Liang Bai, Chen Li, and Xian Yang. A deep clustering-based  
247 state-space model for improved disease risk prediction in personalized healthcare. *Annals of  
248 Operations Research*, 341(1):647–672, 2024.
- 249 [7] Daniel Cunha Oliveira, Dylan Sandfelder, André Fujita, Xiaowen Dong, and Mihai Cu-  
250 curingu. Tactical asset allocation with macroeconomic regime detection. *arXiv preprint  
251 arXiv:2503.11499*, 2025.
- 252 [8] Rongbo Chen, Mingxuan Sun, Kunpeng Xu, Jean-Marc Patenaude, and Shengrui Wang.  
253 Clustering-based cross-sectional regime identification for financial market forecasting. In  
254 *International Conference on Database and Expert Systems Applications*, pages 3–16. Springer,  
255 2022.
- 256 [9] Jose MG Vilar and Leonor Saiz. Dynamics-informed deconvolutional neural networks for super-  
257 resolution identification of regime changes in epidemiological time series. *Science Advances*, 9  
258 (28):eadf0673, 2023.
- 259 [10] Carolyn Augusta, Rob Deardon, and Graham Taylor. Deep learning for supervised classification  
260 of spatial epidemics. *Spatial and Spatio-temporal Epidemiology*, 29:187–198, 2019.
- 261 [11] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, Yoshua  
262 Bengio, et al. Graph attention networks. *stat*, 1050(20):10–48550, 2017.
- 263 [12] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
264 networks? In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryGs6iA5Km>.
- 265
- 266 [13] Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kaneza-  
267 shi, Tim Kaler, Tao Schardl, and Charles Leiserson. EvolveGCN: Evolving graph convolutional  
268 networks for dynamic graphs. *Proc. Conf. AAAI Artif. Intell.*, 34(04):5363–5370, April 2020.
- 269 [14] Jiaxuan You, Tianyu Du, and Jure Leskovec. ROLAND: Graph learning framework for dynamic  
270 graphs. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and  
271 Data Mining*, pages 2358–2366, New York, NY, USA, August 2022. ACM.
- 272 [15] Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured  
273 state spaces. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=uYLFoz1v1AC>.
- 274
- 275 [16] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In  
276 *International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=AL1fq05o7H>.
- 277

- 278 [17] Gongpei Zhao, Tao Wang, Yi Jin, Congyan Lang, Yidong Li, and Haibin Ling. Grassnet: State  
279 space model meets graph neural network. *CoRR*, abs/2408.08583, 2024. doi: 10.48550/ARXIV.  
280 2408.08583. URL <https://doi.org/10.48550/arXiv.2408.08583>.
- 281 [18] Ali Behrouz and Farnoosh Hashemi. Graph mamba: Towards learning on graphs with state  
282 space models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery*  
283 *and Data Mining*, volume 24, pages 119–130, New York, NY, USA, August 2024. ACM.
- 284 [19] Lincan Li, Hanchen Wang, Wenjie Zhang, and Adelle Coster. Stg-mamba: Spatial-temporal  
285 graph learning via selective state space model. *arXiv preprint arXiv:2403.12418*, 2024.
- 286 [20] Haonan Yuan, Qingyun Sun, Zhaonan Wang, Xingcheng Fu, Cheng Ji, Yongjian Wang, Bo Jin,  
287 and Jianxin Li. DG-Mamba: Robust and efficient dynamic graph structure learning with  
288 selective state space models. *Proc. Conf. AAAI Artif. Intell.*, 39(21):22272–22280, April 2025.
- 289 [21] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax.  
290 In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France,*  
291 *April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL [https://](https://openreview.net/forum?id=rkE3y85ee)  
292 [openreview.net/forum?id=rkE3y85ee](https://openreview.net/forum?id=rkE3y85ee).
- 293 [22] Junyoung Chung, Çağlar Gülçehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation  
294 of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014. URL  
295 <http://arxiv.org/abs/1412.3555>.
- 296 [23] Saeed Masoudnia and Reza Ebrahimpour. Mixture of experts: a literature survey. *Artificial*  
297 *Intelligence Review*, 42(2):275–293, 2014.
- 298 [24] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio  
299 and Yann LeCun, editors, *2nd International Conference on Learning Representations, ICLR*  
300 *2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. URL  
301 <http://arxiv.org/abs/1312.6114>.
- 302 [25] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov,  
303 and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30,  
304 2017.

## 305 5 Appendix

### 306 5.1 Training Procedure Details

307 **Decoder routing and teacher forcing.** During training the MoE is routed by a convex state mix

$$s_t = (1 - \eta) p_\theta(z_t | h_G^{t-1}, z_{t-1}) + \eta q_\phi(z_t | h_G^t), \quad (2)$$

308 with  $\eta \in [0, 1]$ . Early in training the decoder relies more on the posterior (teacher forcing,  $\eta \approx 1/3$ );  
309 as training progresses,  $\eta$  is annealed to 0 so emissions are governed by the learned prior, matching  
310 the causal rollout used at inference. This reduces exposure bias without sacrificing stability.

311 **Temperature and alignment schedules.** We parameterize  $q_\phi$  with a Gumbel–Softmax at temper-  
312 ature  $\tau$ . We anneal  $\tau$  from 1.0 to a small floor (e.g., 0.2) over the first half of training to promote  
313 confident, non-degenerate state usage while avoiding premature hard assignments. The alignment  
314 weight  $\beta$  is linearly warmed from 0 to 1 over the first third of training so that forecasting stabilizes  
315 before the encoder–prior terms dominate. The prior-fitting weight  $\gamma$  is set to 1 and may be mildly  
316 reduced later (e.g., to 0.7 after 60% of training) if the learned prior becomes too reactive. The decoder  
317 mix  $\eta$  is annealed linearly to 0 over the first 40% of training to phase out teacher forcing. These  
318 schedules were chosen empirically to prevent posterior collapse, avoid chasing a moving prior, and  
319 align the training-time routing with the inference-time strictly time-ordered rollout.

320 **Lightweight regularization.** We add two small regularizers that do not alter the loss but improve  
321 state usage: (i) a *diversity* term that keeps the batch-average posterior close to uniform,  $\text{KL}(\bar{q} \| \text{Unif})$   
322 with  $\bar{q} = \frac{1}{B} \sum_i q_\phi^{(i)}(z_t | h_G^t)$ , to avoid dead states; and (ii) a *sharpness* term that lowers the entropy  
323 of per-graph posteriors,  $\mathbb{E}[H(q_\phi(z_t | h_G^t))]$ , ramped in after the temperature has decreased. Both are  
324 coefficients of small amplitudes (e.g.,  $\lambda_{\text{marg}} \approx 0.1$ ,  $\lambda_{\text{sharp}} \leq 0.05$ ).

325 **Implementation notes.** The temporal encoder is instantiated with a ROLAND-style [14] dynamic  
326 GNN that maintains node memories via GRU updates and pools to  $h_G^t$ , but we discard the the live  
327 update and caching mechanisms that are not relevant in our setup. All training were performed using  
328 internal cluster GPUs. A couple of workers (2-4) are sufficient due to the small size of the dataset.  
329 The dataset was randomly split with label stratification (based on regime) following a 75%/25% split  
330 for training/validation. Hyperparameters of each method were selected using 4-fold cross-validation  
331 on the training set.

332 The MoE decoder comprises  $K$  independent MLP experts  $\{f_k\}_{k=1}^K$  with outputs combined by  $s_t$   
333 from Eq. (2). Experts can be replaced with domain-specific modules without changing the objective.  
334 We optimize with Adam, apply gradient clipping for stability, and select checkpoints on validation  
335 one-step error.

### 336 5.2 Dataset Generation Process

337 We generate 150 spatio-temporal graph sequences with  $T = 10$  snapshots (TP1–TP10). Each  
338 sequence begins with a randomly sampled discrete regime  $r_1 \in \{r_1, r_2, r_3\}$ , and undergoes a single  
339 potential regime change before TP5, as detailed below.

340 **Node clusters.** Let  $d = 8$  denote the feature dimension. We sample the number of clusters  
341  $C \sim \text{Unif}\{3, 4, 5\}$ . For each cluster  $c \in \{1, \dots, C\}$ , we draw a center  $\mu_c \sim \text{Unif}([-10, 10]^d)$  and  
342 use an isotropic covariance  $\Sigma = 0.6^2 I_d$ . We then sample the cluster size  $M_c \sim \text{Unif}\{5, \dots, 100\}$   
343 and node coordinates

$$x_i^{(1)} \sim \mathcal{N}(\mu_{c(i)}, \Sigma) \quad \text{for } i = 1, \dots, \sum_{c=1}^C M_c.$$

344 **Edges (spatial proximity).** For each snapshot  $t$ , we build an undirected  $k$ -nearest neighbor graph  
345 on  $\{x_i^{(t)}\}_i$  in  $\mathbb{R}^d$  with Euclidean distance and

$$k = \max_{1 \leq c \leq C} M_c + 1$$

346 to assure bridges between clusters. We keep the same set of edges the graph at each  $t$ ; since the  
 347 dynamics below are global translations plus small noise, the topology is translation-invariant and  
 348 empirically stable across  $t$ .

349 **Regimes and dynamics.** Regimes induce constant drifts along all coordinates:

$$v(r_1) = -2 \mathbf{1}_d, \quad v(r_2) = +2 \mathbf{1}_d, \quad v(r_3) = \mathbf{0}_d.$$

350 Let  $\varepsilon_i^{(t)} \sim \mathcal{N}(0, 0.6^2 I_d)$  be i.i.d. perturbations. For  $t = 1, \dots, 9$ , node positions evolve as

$$x_i^{(t+1)} = x_i^{(t)} + v(r_t) + \varepsilon_i^{(t)}.$$

351 Before applying the update to obtain TP5 (i.e., between TP4 and TP5), we resample the regime  
 352  $r_k^{TP5} \sim \text{Unif}\{r_1, r_2, r_3\}$  independently of  $r_k^{TP1}$ ; consequently, some sequences keep their regime  
 353 while others switch. The labels  $\{r_k^{TPi}\}$  are not used for supervision.

354 The dataset was randomly split with label stratification (based on regime) following a 75%/25% split  
 355 for training/validation.

356 **Rationale.** This construction test models’ ability to infer discrete regime trajectories from spatially  
 357 structured observations and to remain robust to a mid-sequence non-stationarity (TP4→TP5).

### 358 5.3 Baselines Implementation Details

359 **MLP (no space, no time)** The MLP receives a concatenation of node features across the observed  
 360 horizon, so the per-node input dimensionality grows linearly with time (e.g., with base feature size  
 361  $d = 8$ , inputs are 8 at TP1, 16 at TP2, etc.). To accommodate dynamic topology (variable node  
 362 counts across graphs and time), we mimic message passing with *self-loops only*: a shared per-node  
 363 MLP processes each node independently (no neighbor aggregation), producing per-node embeddings  
 364 at time  $t$ . We then apply parametric pooling via a small pooling MLP (DeepSets-style [25]) to obtain  
 365 a graph-level context vector. Final node-level predictions  $\hat{X}_{t+1}$  are produced by another shared MLP  
 366 that conditions on both the node’s self-updated embedding and the pooled context. This design  
 367 ignores explicit topology while still permitting information mixing through learnable pooling.

368 **Persistence** ( $X_{t+1} = X_t$ ) A parameter-free, time-only baseline that copies the last observation to  
 369 the next step. It has a slight advantage in regimes with near-constant dynamics ( $r_3$  in the synthetic  
 370 dataset, where the drift is absent) but remains weak overall, providing a lower bound that tests whether  
 371 temporal autocorrelation alone explains performance.

372 **LSTM-GNN (ROLAND-based)** We instantiate a ROLAND-style dynamic GNN encoder ([14])  
 373 with GRU updates ([22]) for hierarchical node states. At each time step, node embeddings are  
 374 updated by a graph layer and then temporally evolved via GRUs; the model natively supports  
 375 dynamic topology. As in our main architecture, we discard caching and live-update mechanisms. For  
 376 this baseline we directly project to node features with a linear head (no graph-level pooling), yielding  
 377 a strong spatio-temporal encoder without discrete regimes or state-conditioned emissions.

378 **STG-Mamba** We follow STG-Mamba [19]: blocks interleave spatial mixing (graph filter-  
 379 ing/propagation on the current adjacency) with temporal Mamba modules that implement selective  
 380 state-space updates along time. Each block uses residual connections, normalization, and pointwise  
 381 MLPs. Stacking several blocks yields multi-scale spatio-temporal modeling. Training minimizes  
 382 next-step MSE. Other GNN-SSM hybrids were considered, but most lacked robustness to topological  
 383 change (relevant for our real-world, public results not yet available) or had no public implementations,  
 384 so we did not include them.

385 All baselines follow a similar training procedure as described in section 5.1 and 2.3. Notably all  
 386 trainings were performed on internal cluster with GPUs. The hyperparameters of each baseline were  
 387 selected using 4-fold cross-validation on the training set, later evaluated on the validation set as  
 388 reported in table 1.

## 389 **NeurIPS Paper Checklist**

### 390 **1. Claims**

391 Question: Do the main claims made in the abstract and introduction accurately reflect the  
392 paper's contributions and scope?

393 Answer: [\[Yes\]](#)

394 Justification: We claim that current methods often disregards regime detection and tracking in  
395 spatio-temporal systems. We show how our method leverage the regimes trajectory through  
396 the discrete state space in Section 3.4, and how it does not interfere with performances in  
397 section 3.3.

398 Guidelines:

- 399 • The answer NA means that the abstract and introduction do not include the claims  
400 made in the paper.
- 401 • The abstract and/or introduction should clearly state the claims made, including the  
402 contributions made in the paper and important assumptions and limitations. A No or  
403 NA answer to this question will not be perceived well by the reviewers.
- 404 • The claims made should match theoretical and experimental results, and reflect how  
405 much the results can be expected to generalize to other settings.
- 406 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
407 are not attained by the paper.

### 408 **2. Limitations**

409 Question: Does the paper discuss the limitations of the work performed by the authors?

410 Answer: [\[Yes\]](#)

411 Justification: Yes in the end of the conclusion section 4 we give some drawbacks and  
412 limitations of the methods together with some possible improvements.

413 Guidelines:

- 414 • The answer NA means that the paper has no limitation while the answer No means that  
415 the paper has limitations, but those are not discussed in the paper.
- 416 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 417 • The paper should point out any strong assumptions and how robust the results are to  
418 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
419 model well-specification, asymptotic approximations only holding locally). The authors  
420 should reflect on how these assumptions might be violated in practice and what the  
421 implications would be.
- 422 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
423 only tested on a few datasets or with a few runs. In general, empirical results often  
424 depend on implicit assumptions, which should be articulated.
- 425 • The authors should reflect on the factors that influence the performance of the approach.  
426 For example, a facial recognition algorithm may perform poorly when image resolution  
427 is low or images are taken in low lighting. Or a speech-to-text system might not be  
428 used reliably to provide closed captions for online lectures because it fails to handle  
429 technical jargon.
- 430 • The authors should discuss the computational efficiency of the proposed algorithms  
431 and how they scale with dataset size.
- 432 • If applicable, the authors should discuss possible limitations of their approach to  
433 address problems of privacy and fairness.
- 434 • While the authors might fear that complete honesty about limitations might be used by  
435 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
436 limitations that aren't acknowledged in the paper. The authors should use their best  
437 judgment and recognize that individual actions in favor of transparency play an impor-  
438 tant role in developing norms that preserve the integrity of the community. Reviewers  
439 will be specifically instructed to not penalize honesty concerning limitations.

### 440 **3. Theory assumptions and proofs**

441 Question: For each theoretical result, does the paper provide the full set of assumptions and  
442 a complete (and correct) proof?

443 Answer: [NA]

444 Justification: No theoretical theorem or results provided in the paper.

445 Guidelines:

- 446 • The answer NA means that the paper does not include theoretical results.
- 447 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
448 referenced.
- 449 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 450 • The proofs can either appear in the main paper or the supplemental material, but if  
451 they appear in the supplemental material, the authors are encouraged to provide a short  
452 proof sketch to provide intuition.
- 453 • Inversely, any informal proof provided in the core of the paper should be complemented  
454 by formal proofs provided in appendix or supplemental material.
- 455 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 456 4. Experimental result reproducibility

457 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
458 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
459 of the paper (regardless of whether the code and data are provided or not)?

460 Answer: [Yes]

461 Justification: We extensively described the generation of the dataset together with imple-  
462 mentation details of both our method and baselines in the Sections 3.1, 2.2, 3.2. We provide  
463 additional information on training procedure, scheduling and implementation in Section 5.2,  
464 5.1, 5.3.

465 Guidelines:

- 466 • The answer NA means that the paper does not include experiments.
- 467 • If the paper includes experiments, a No answer to this question will not be perceived  
468 well by the reviewers: Making the paper reproducible is important, regardless of  
469 whether the code and data are provided or not.
- 470 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
471 to make their results reproducible or verifiable.
- 472 • Depending on the contribution, reproducibility can be accomplished in various ways.  
473 For example, if the contribution is a novel architecture, describing the architecture fully  
474 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
475 be necessary to either make it possible for others to replicate the model with the same  
476 dataset, or provide access to the model. In general, releasing code and data is often  
477 one good way to accomplish this, but reproducibility can also be provided via detailed  
478 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
479 of a large language model), releasing of a model checkpoint, or other means that are  
480 appropriate to the research performed.
- 481 • While NeurIPS does not require releasing code, the conference does require all submis-  
482 sions to provide some reasonable avenue for reproducibility, which may depend on the  
483 nature of the contribution. For example
  - 484 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
485 to reproduce that algorithm.
  - 486 (b) If the contribution is primarily a new model architecture, the paper should describe  
487 the architecture clearly and fully.
  - 488 (c) If the contribution is a new model (e.g., a large language model), then there should  
489 either be a way to access this model for reproducing the results or a way to reproduce  
490 the model (e.g., with an open-source dataset or instructions for how to construct  
491 the dataset).
  - 492 (d) We recognize that reproducibility may be tricky in some cases, in which case  
493 authors are welcome to describe the particular way they provide for reproducibility.  
494 In the case of closed-source models, it may be that access to the model is limited in

495 some way (e.g., to registered users), but it should be possible for other researchers  
496 to have some path to reproducing or verifying the results.

#### 497 **5. Open access to data and code**

498 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
499 tions to faithfully reproduce the main experimental results, as described in supplemental  
500 material?

501 Answer: [No]

502 Justification: The process for generating data is explained in extensive details, allowing  
503 replication. The code base is kept private until final publication.

504 Guidelines:

- 505 • The answer NA means that paper does not include experiments requiring code.
- 506 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
507 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 508 • While we encourage the release of code and data, we understand that this might not be  
509 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
510 including code, unless this is central to the contribution (e.g., for a new open-source  
511 benchmark).
- 512 • The instructions should contain the exact command and environment needed to run to  
513 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
514 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 515 • The authors should provide instructions on data access and preparation, including how  
516 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 517 • The authors should provide scripts to reproduce all experimental results for the new  
518 proposed method and baselines. If only a subset of experiments are reproducible, they  
519 should state which ones are omitted from the script and why.
- 520 • At submission time, to preserve anonymity, the authors should release anonymized  
521 versions (if applicable).
- 522 • Providing as much information as possible in supplemental material (appended to the  
523 paper) is recommended, but including URLs to data and code is permitted.

#### 524 **6. Experimental setting/details**

525 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
526 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
527 results?

528 Answer: [Yes]

529 Justification: All the details are present in the appendix (Section 5.1, 5.2, 5.3). We keep a  
530 couple of specific implementation details and hyperparameters search for the main publica-  
531 tion.

532 Guidelines:

- 533 • The answer NA means that the paper does not include experiments.
- 534 • The experimental setting should be presented in the core of the paper to a level of detail  
535 that is necessary to appreciate the results and make sense of them.
- 536 • The full details can be provided either with the code, in appendix, or as supplemental  
537 material.

#### 538 **7. Experiment statistical significance**

539 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
540 information about the statistical significance of the experiments?

541 Answer: [NA]

542 Justification: The set of experiments reported in this paper is directly performed on the  
543 left-out test set after cross-validation on the training set. Thus no direct measure of statistical  
544 significance can be performed on the results reported in Table 1. We could report the detailed  
545 cross-validation results if the reviewers wish it.

546 Guidelines:

- 547 • The answer NA means that the paper does not include experiments.
- 548 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 549 dence intervals, or statistical significance tests, at least for the experiments that support
- 550 the main claims of the paper.
- 551 • The factors of variability that the error bars are capturing should be clearly stated (for
- 552 example, train/test split, initialization, random drawing of some parameter, or overall
- 553 run with given experimental conditions).
- 554 • The method for calculating the error bars should be explained (closed form formula,
- 555 call to a library function, bootstrap, etc.)
- 556 • The assumptions made should be given (e.g., Normally distributed errors).
- 557 • It should be clear whether the error bar is the standard deviation or the standard error
- 558 of the mean.
- 559 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 560 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 561 of Normality of errors is not verified.
- 562 • For asymmetric distributions, the authors should be careful not to show in tables or
- 563 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 564 error rates).
- 565 • If error bars are reported in tables or plots, The authors should explain in the text how
- 566 they were calculated and reference the corresponding figures or tables in the text.

## 567 8. Experiments compute resources

568 Question: For each experiment, does the paper provide sufficient information on the com-  
 569 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
 570 the experiments?

571 Answer: [Yes]

572 Justification: In the appendix (5.1), we provide information about annealing rate and general  
 573 details on the training procedure.

574 Guidelines:

- 575 • The answer NA means that the paper does not include experiments.
- 576 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 577 or cloud provider, including relevant memory and storage.
- 578 • The paper should provide the amount of compute required for each of the individual
- 579 experimental runs as well as estimate the total compute.
- 580 • The paper should disclose whether the full research project required more compute
- 581 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 582 didn't make it into the paper).

## 583 9. Code of ethics

584 Question: Does the research conducted in the paper conform, in every respect, with the  
 585 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

586 Answer: [Yes]

587 Justification: We acknowledge the NeurIPS Code of Ethics and believe that the method  
 588 described in this paper has no specific negative societal impact and potential harmful  
 589 consequences.

590 Guidelines:

- 591 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 592 • If the authors answer No, they should explain the special circumstances that require a
- 593 deviation from the Code of Ethics.
- 594 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 595 eration due to laws or regulations in their jurisdiction).

## 596 10. Broader impacts

597 Question: Does the paper discuss both potential positive societal impacts and negative  
 598 societal impacts of the work performed?

599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651

Answer: [Yes]

Justification: We discussed how our work can be helpful in high-stake applications (Sections 1 and 4) where the interpretability of the spatio-temporal dynamics can be mapped to well established regimes, like in healthcare or finance.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risk.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We credited all previous work, especially previous published methods, that were used throughout the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.

- 652 • The authors should cite the original paper that produced the code package or dataset.
- 653 • The authors should state which version of the asset is used and, if possible, include a
- 654 URL.
- 655 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 656 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 657 service of that source should be provided.
- 658 • If assets are released, the license, copyright information, and terms of use in the
- 659 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)
- 660 has curated licenses for some datasets. Their licensing guide can help determine the
- 661 license of a dataset.
- 662 • For existing datasets that are re-packaged, both the original license and the license of
- 663 the derived asset (if it has changed) should be provided.
- 664 • If this information is not available online, the authors are encouraged to reach out to
- 665 the asset’s creators.

### 666 13. New assets

667 Question: Are new assets introduced in the paper well documented and is the documentation  
668 provided alongside the assets?

669 Answer: [Yes]

670 Justification: The generated synthetic dataset has been extensively described in section 5.2.

671 Guidelines:

- 672 • The answer NA means that the paper does not release new assets.
- 673 • Researchers should communicate the details of the dataset/code/model as part of their
- 674 submissions via structured templates. This includes details about training, license,
- 675 limitations, etc.
- 676 • The paper should discuss whether and how consent was obtained from people whose
- 677 asset is used.
- 678 • At submission time, remember to anonymize your assets (if applicable). You can either
- 679 create an anonymized URL or include an anonymized zip file.

### 680 14. Crowdsourcing and research with human subjects

681 Question: For crowdsourcing experiments and research with human subjects, does the paper  
682 include the full text of instructions given to participants and screenshots, if applicable, as  
683 well as details about compensation (if any)?

684 Answer: [NA]

685 Justification: No involvement of crowdsourcing or research with human subjects.

686 Guidelines:

- 687 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 688 human subjects.
- 689 • Including this information in the supplemental material is fine, but if the main contribu-
- 690 tion of the paper involves human subjects, then as much detail as possible should be
- 691 included in the main paper.
- 692 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 693 or other labor should be paid at least the minimum wage in the country of the data
- 694 collector.

### 695 15. Institutional review board (IRB) approvals or equivalent for research with human 696 subjects

697 Question: Does the paper describe potential risks incurred by study participants, whether  
698 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
699 approvals (or an equivalent approval/review based on the requirements of your country or  
700 institution) were obtained?

701 Answer: [NA]

702 Justification: No involvement of crowdsourcing or research with human subjects.

703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725

**Guidelines:**

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

**16. Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development is not centered around LLMs.

**Guidelines:**

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.