

STRUCTURAL KNOWLEDGE INFORMED CONTINUAL MULTIVARIATE TIME SERIES FORECASTING

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent studies in multivariate time series (MTS) forecasting reveal that explicitly modeling the hidden dependencies among different time series can yield promising forecasting performance and reliable explanations. However, modeling variable dependencies remains underexplored when a sequence of MTS under different regimes (stages) is continuously accumulated. Due to the potential distribution and dependency disparities, the underlying model may encounter the catastrophic forgetting problem, *i.e.*, it is challenging to memorize and infer different types of variable dependencies across different regimes while maintaining forecasting performance. To this end, we propose a novel Structural Knowledge Informed Continual Learning (SKI-CL) framework to perform MTS forecasting under the continual learning setting, which leverages the structural knowledge to characterize the dynamic variable dependencies within each regime. Specifically, we first develop a deep forecasting model with a graph learner that enables fine-grained dynamic MTS modeling. Next, we impose a regularization scheme to ensure the consistency between the learned variable dependencies and the structure knowledge (*e.g.*, physical constraints, domain knowledge, feature similarity). Finally, we develop a representation-matching memory replay scheme to tackle the catastrophic forgetting problem, which maximizes the temporal coverage of MTS data to efficiently preserve the underlying temporal dynamics and dependency structures of each regime. Thoroughly empirical studies on synthetic and real-world benchmarks validate SKI-CL’s efficacy and advantages over the state-of-the-art in tackling continual MTS forecasting tasks. In addition, SKI-CL can accurately infer learned dependencies in the test stage based on MTS data without knowing the identities of different regimes.

1 INTRODUCTION

Multivariate time series (MTS) forecasting aims to predict the future horizons of multiple time series based on their historical values and has shown its importance in various applications, *e.g.*, healthcare, traffic control, energy management, and finance Jin et al. (2018); Gonzalez-Vidal et al. (2019); Guo et al. (2019); Zhang et al. (2017). Accurate forecasting not only relies on capturing temporal dynamics of the historical multivariate time series data, *e.g.*, autoregressive integrated moving average (ARIMA) Box et al. (2015), deep convolutional neural networks van den Oord et al. (2016); Bai et al. (2018); Borovykh et al. (2017), recurrent neural networks Bai et al. (2018); Lai et al. (2018); Shih et al. (2019), and transformer models Zhou et al. (2021); Wu et al. (2021); Zhou et al. (2022); Liu et al. (2021), but also relies on modeling dependency structures among different time series Graves (2013); Lai et al. (2018); Shih et al. (2019). However, the nonstationary nature of most MTS data implies that the structural dependencies among different time series vary across time. To make it worse, in real-world applications, different regimes (stages) of MTS data are often continuously collected under different internal logic of the target system.

To develop an effective forecasting model that can handle continuously accumulated MTS data from different regimes and infer the dependency structures accordingly, there are two fundamental challenges. The first one is how to model dependency structures among different time series by incorporating the structural knowledge of the target system. Recently, there is increasing evidence suggesting that jointly optimizing the inferred structure toward the existing structural knowledge and

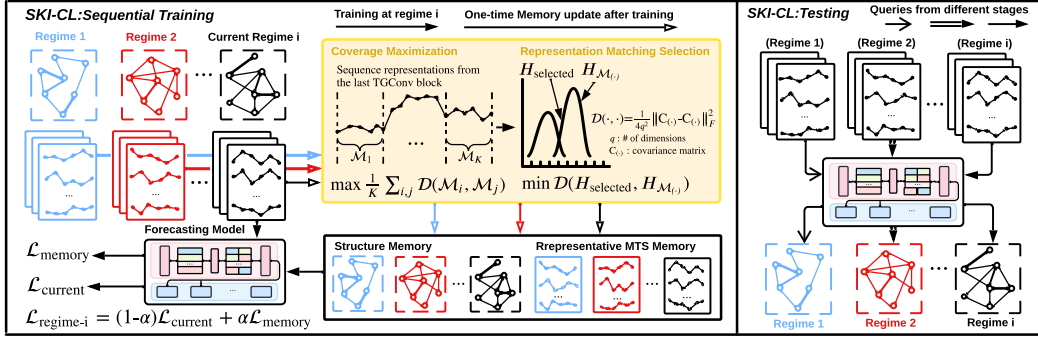


Figure 1: The proposed SKI-CL framework for continual MTS forecasting. The training objectives for each regime contains the current training data and the memory buffer. After training at each regime, the structural knowledge and samples selected by our representation-matching scheme are added to the current memory. At testing phase, SKI-CL is able to dynamically infer faithful dependency structures for different regimes without accessing the memory buffer.

forecasting objectives can render reliable explanations and competing performance Shang & Chen (2021). The structural knowledge can be available in different formats, *e.g.*, physical constraints such as traffic networks, power grids, and sensor networks Li et al. (2017); Luo et al. (2021); Khodayar & Wang (2018); Yan et al. (2018), or dependencies (correlations) that are inferred or derived from raw data by leveraging either domain knowledge or traditional statistical methods Chen et al. (2022); Duan et al. (2022); Lin et al. (2021); Cao et al. (2020); Shang & Chen (2021). The second challenge is how to keep track of the latest regime of MTS data while maintaining the performance on existing regimes. One potential solution is to retrain the forecasting model periodically over the newly collected regime to incorporate the latest structural knowledge. This method, however, will lead to the catastrophic forgetting problem, *i.e.*, the model performance will deteriorate over existing regimes as their associated structural knowledge cannot be maintained. Note that joint training may not be a feasible solution due to the memory constraints of storing all historical data (different regimes) and the computational complexity of handling an ever-increasing number of diverse scenarios.

In this paper, we present a novel Structural Knowledge Informed Continual Learning (SKI-CL) framework that incrementally learns and preserves meaningful dependency structures for MTS forecasting. As shown in Figure 1, the key idea is to exploit structural knowledge to characterize the variable dependencies within each different regime. Specifically, we first design a dynamic graph structure learning module that captures the temporal patterns to construct the parameterized interactions between all variable pairs. Instead of learning a single fixed graph structure Bai et al. (2020); Wu et al. (2020), our module dynamically infers dependency structures (in the form of graphs) based on different MTS input windows to cope with the dependencies variations within each regime, thus enabling non-stationary MTS modeling. We further enforce the consistency between the learned structures and the structural knowledge by designing an adaptive regularization scheme that accommodates different forms of structural knowledge that exist in MTS. Four common scenarios regarding the description of dependencies and the availability of structural knowledge is considered, *i.e.*, the discrete/continuous edge description, and fully/partially observed structural knowledge. Finally, we present a novel representation-matching memory replay scheme to tackle the catastrophic forgetting problem, which maximizes the temporal coverage of MTS data to efficiently preserve the underlying temporal dynamics and dependency structure of each regime (Figure 1(middle)). By jointly learning from the current regime and the constructed memory, the obtained model is able to maintain the forecasting performance and infer the learned structures from the existing regimes.

In summary, our work makes the following four contributions. (1) We present a novel Structural Knowledge Informed Continual Learning (SKI-CL) framework to perform MTS forecasting and infer dependency structures under the continual learning setting. (2) We develop a dynamic graph structure learning module to capture temporal dependencies and infer dependency structures, and employ an adaptive regularization scheme to incorporate structural knowledge. (3) We propose a novel representation-matching memory replay scheme to maximize the temporal coverage of MTS data and preserve the underlying temporal dynamics and dependency structures of each regime. (4) Thoroughly experiments on one synthetic dataset and three benchmark datasets demonstrate the superiority of SKI-CL over the state-of-the-art in continual MTS forecasting and dependency structure inference.

2 RELATED WORK

2.1 MODELING DEPENDENCIES IN MULTIVARIATE TIME SERIES FORECASTING

In recent years, modeling variable dependencies of MTS has received increasing attention for forecasting tasks. Early methods apply linear or convolution transformations to capture variable dependencies in an implicit recurrent process Lai et al. (2018); Graves (2013); Shih et al. (2019), which fall short of modeling the non-Euclidean interactions due to the underlying fully-connected or translation-invariant assumptions. The advent of Graph Neural Networks (GNNs) has inspired the formulation of variable dependencies as a given or learnable graph, with variables being nodes and pairwise relationships being edges. Existing literature models the dependency structures based on different topological perspectives and temporal granularity (*e.g.*, undirected Yu et al. (2018) and directed graph Li et al. (2017), static Bai et al. (2020); Wu et al. (2020) and dynamic graphs Ye et al. (2022); Cao et al. (2020), single or multiple layers Lin et al. (2021)). On the other hand, structural knowledge has been an important component in GNN-based forecasting methods. In many tasks such as the traffic Guo et al. (2019); Li et al. (2017) and skeleton-based action prediction Yan et al. (2018), structural knowledge is explicitly presented as spatial connections. In other cases without an explicit topological structure, the structural knowledge can be drawn from either domain knowledge (*e.g.*, transfer entropy Duan et al. (2022), Mel-frequency cepstral coefficients Lin et al. (2021)) or feature similarity (*e.g.*, the correlations of decomposed time series Ng et al. (2022), kNN graph Shang & Chen (2021)). Their promising results suggest the capability of structural knowledge to convey meaningful dependency information. Our proposed method enforces the consistency between the learned graph structures and the structural knowledge so as to characterize the underlying relation-temporal dependencies and improve the continual MTS forecasting performance.

2.2 CONTINUAL LEARNING IN MULTIVARIATE TIME SERIES FORECASTING

Deep learning models resort to continual learning to address the catastrophic forgetting issue when sequentially adapting to new tasks. Existing literature in continual learning lies in three folds: the experience-replayed methods Rolnick et al. (2019), the parameter-isolation methods Rusu et al., and regularization-based methods Kirkpatrick et al. (2017); Li & Hoiem (2017). Current continual learning works have been extensively studied on images Wang et al. (2022), texts Ke & Liu (2022), and graph data Zhang et al. (2022a;b). However, much less attention is drawn to time series data, and the focus has primarily been on classification and forecasting tasks without explicitly addressing complex variable dependencies Gupta et al. (2021); He & Sick (2021). How to maintain the meaningful dependency structures and forecasting performance over different regimes is underexplored. Our proposed SKI-CL tackles this issue by jointly optimizing the inferred structure toward the existing structural knowledge and forecasting objectives based on samples drawn from the current regime and a representative memory.

We also notice that recent efforts have been made to perform forecasting Pham et al. (2022); Zhang et al. (2023) and/or topology identification Money et al. (2021); Natali et al. (2022); Isufi et al. (2019); Zaman et al. (2020) from MTS data in an online learning setting, which focus on adapting the forecasting model and dependency structure from the historical MTS data to future unseen data. On the contrary, our study aims to maintain the forecasting performance and infer the learned structures from the existing regimes while continuously updating the model over the latest regime. Without forgetting the structure knowledge that has been acquired previously, the model can easily cope with similar regimes that may be encountered in the future. [In addition to online MTS forecasting, we discuss the relationships between our study and other related research topics in Appendix B.](#)

3 METHODOLOGY

In this section, we first formulate the continual MTS forecasting task with dependency structure learning. Then, we introduce the proposed framework including the structural knowledge-informed graph learning model and the representation-matching sample selection scheme for continual MTS forecasting, as shown in Figure 2 and Figure 1(middle), respectively.

3.1 PROBLEM STATEMENT

We first introduce the MTS forecasting task in a single regime (stage). Let $X \in \mathbb{R}^{N \times T}$ denote the MTS data containing N variables and T total time steps, where $X_{:,t} \in \mathbb{R}^{N \times 1}$ denotes t -th time

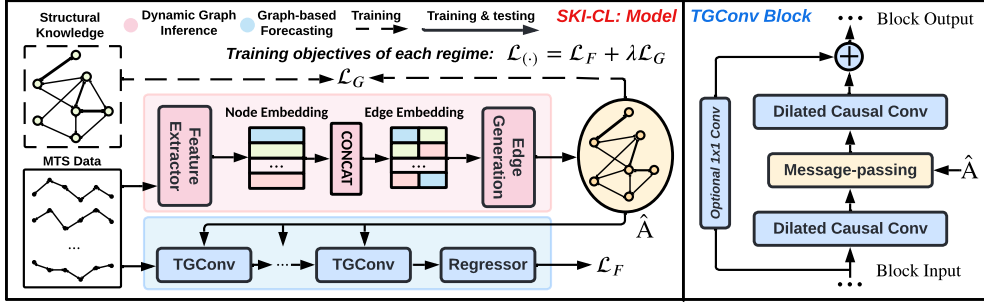


Figure 2: The proposed SKI-CL Model for dependencies modeling and MTS forecasting.

step across all variables and $X_{i,:} \in \mathbb{R}^{1 \times T}$ denotes i -th variable. Our target is to learn a model that includes a dynamic graph inference module $\mathcal{G}(\cdot)$ summarizing a historical τ -step window of MTS as a graph to encode the dependency structure, as well as a forecasting module $\mathcal{F}(\cdot)$ predicting the next τ' time steps based on the input window and inferred graph. Mathematically, at a starting time step t , the corresponding forecast is defined as: $\hat{X}_{:,t+\tau:t+\tau+\tau'-1} = \mathcal{F}(X_{:,t:t+\tau-1}, \mathcal{G}(X_{:,t:t+\tau-1}))$.

We further extend the forecasting task to a continual learning setting. In continual learning, there exist S distinct regimes of MTS data with different dependencies and temporal dynamics. The model can only access MTS data of the current regime. Denoting the data of s -th regime as $X^{(s)}$, the objective is to learn a model to minimize the forecasting error across all seen regimes: $\mathcal{F}^*, \mathcal{G}^* = \arg \min_{\mathcal{F}, \mathcal{G}} \sum_{s=1}^S L(\hat{X}_{:,t:t+\tau-1}^{(s)}, X_{:,t+\tau:t+\tau+\tau'-1}^{(s)})$ with L being the loss function. We assume there is a readily available or extracted structural knowledge $A \in \mathbb{R}^{N \times N}$ (either partial or completed) at each regime that serves as a reference to characterize the underlying dependencies.

3.2 STRUCTURAL KNOWLEDGE INFORMED GRAPH LEARNING FOR MTS FORECASTING

Dynamic Graph Inference Module. Different from the existing works that generate a static graph at the regime level Wu et al. (2020); Bai et al. (2020); Shang & Chen (2021), we aim to model the non-stationary variable dependencies of MTS as a dynamic graph at the granularity of an input window. Therefore, as shown in Figure 2, we construct a dynamic graph inference module that more precisely reveals the relation-temporal dynamics in a single regime and has the capacity to handle dependencies change when the regime shifts.

In this module, each edge is explicitly modeled for all node pairs, which is parameterized as follows. For a possible edge connecting node i and j , we use the stacked convolution p times along the temporal dimension as a feature extractor to yield node embedding z_i and z_j (i.e., $z_i = (\mathbf{FC} \circ \mathbf{Conv}^{[p]})(X_{i,t:t+\tau-1})$) which are concatenated as the edge embedding. Next, p' fully connected layers finalize the generation of a parameterized edge as $\hat{A}_{ij} = \mathbf{FC}^{[p']}(z_i || z_j)$. The output graph $\hat{A} \in \mathbb{R}^{N \times N}$ summarizes the variable interactions from the temporal dynamics within a sequence, and can be further utilized to generate the forecasts.

Incorporating Structural Knowledge for Dependencies Characterization. To learn a faithful dynamic graph structure that characterizes the underlying dependencies of each regime, we incorporate structural knowledge as a reference in learning objectives. In real-world MTS modeling, the edges can be either continuous if we can quantify the strength in the context, or binary if we are more confident of the connection in a qualitative sense (e.g., physical connections). Moreover, it also interleaves with the fact that if the structural knowledge can be fully observed, as in many cases we are only confident in the existence of certain relationships.

To fully leverage different forms of structural knowledge in dependency structure learning, we design an adaptive scheme that imposes different constraints on the parameterized graph in the objective function, which is denoted as \mathcal{L}_G as shown in Figure 2 (left). If an edge is treated as a binary variable, we activate the parameterized edge with a sigmoid function to approximate the Bernoulli distribution $\hat{A}_{ij} \sim \text{Bern}(\theta_{ij})$, where $P(\hat{A}_{ij} = 1) = \theta_{ij}$ is the probability that an edge is forming between node i and node j . Then, we encourage the probability of edge to be consistent with the prior, which essentially minimizes the binary cross entropy: $\mathcal{L}_G = \sum_{i,j} -A_{ij} \log \theta_{ij} - (1 - A_{ij}) \log (1 - \theta_{ij})$.

If an edge is treated as a continuous variable, we activate the parameterized edge with a ReLU function to remove the weak connections, and enforce the consistency between the numerical values of the parameterized edge and prior, representing a similar interaction strength. This is achieved by the MSE objective $\mathcal{L}_G = \frac{1}{N^2} \sum_{i,j} \|A_{ij} - \hat{A}_{ij}\|^2$. So far, we have discussed the cases when structural knowledge is readily available/fully observed. For partially observed structural knowledge, we only enforce the consistency between the known entries in the structural knowledge and corresponding parameterized edges, as the dynamic graph inference module is still able to capture and infer the underlying dynamic dependencies via optimizing based on the existing structural knowledge and the forecasting objective.

Graph-based Forecasting Module. To further exploit the structural and temporal dependencies and render forecasts, we design a graph-based forecasting module consisting of multiple Temporal Graph Convolution (TGConv) blocks, as shown in Figure 2. In each block, we leverage a dilated causal convolution Bai et al. (2018); van den Oord et al. (2016) to effectively captures forward dynamics of time series. The dilated causal convolution operation on a 1D sequence input \mathbf{h} is expressed as $r_t = \sum_{k=0}^{K-1} f(k) \cdot h_{t-d-k}$, where r_t denotes the t -th step of obtained representation \mathbf{r} , d represents dilation factor, $f(k)$ is the convolution kernel with size k . Since dilated causal convolution exclusively processes univariate time series, we facilitate the modeling of dependencies in MTS by exchanging and aggregating information through the learned structure, denoted as \hat{A} . This is achieved by a simple yet effective message-passing neural operation Morris et al. (2019), given the collection of univariate representations $(\mathbf{r}_1, \dots, \mathbf{r}_N)$:

$$\text{MessagePassing}_{\hat{A}}(\mathbf{r}_i) = \mathbf{W}_1 \mathbf{r}_i + \mathbf{W}_2 \sum_{j \in \mathcal{N}(i)} e_{j,i} \cdot \mathbf{r}_j \quad (1)$$

where $\mathcal{N}(i)$ represents the neighbors of variable i , $e_{j,i} \in \hat{A}$ denotes an edge weight, $\mathbf{W}_{(\cdot)}$ denotes the learnable weights, and the bias term is omitted for simplicity.

We stack both operations to construct a TGConv block, with an optional 1×1 convolution tackling the possibly different dimensions between the residual input and output. Finally, a fully connected layer serves as a regressor projecting sequence representations onto the forecasts. We adopt the mean squared error between the forecasts and the ground truths as the main learning objective \mathcal{L}_F . The total learning objective function for each regime consists of the forecasting objective and the consistency regularization weighted by a hyperparameter λ .

$$\mathcal{L}_{total} = \mathcal{L}_F + \lambda \mathcal{L}_G = \frac{1}{\tau'} \sum_{t=t+\tau}^{t=t+\tau+\tau'-1} \|\hat{X}_{:,t} - X_{:,t}\|^2 + \lambda \mathcal{L}_G \quad (2)$$

3.3 REPRESENTATION-MATCHING SAMPLE SELECTION FOR CONTINUAL MTS FORECASTING

To tackle the forgetting of variable dependencies and temporal dynamics in sequential training, we store a small subset of MTS samples and the structural knowledge from the previous regimes for memory replay when adapting the model to the current regime. Specifically, we propose an efficient sample selection scheme that maximizes the temporal and dependencies coverage of each regime given a limited memory budget. According to the principle of maximum entropy Guiasu & Shenitzer (1985), we can best represent the underlying knowledge of MTS in each regime with the largest entropy, namely, with the most diverse partitions/modes of relational and temporal patterns. Inspired by this principle and its success in characterizing temporal distribution Du et al. (2021), we perform a distribution characterization by splitting the MTS data to the most diverse modes on the representation space (*i.e.*, the representation of all time-consecutive sequences/samples that encode variable dependencies and temporal dynamics), which is formulated as a constrained optimization problem:

$$\begin{aligned} & \max_{0 < K \leq K_0} \max_{n_1, \dots, n_K} \frac{1}{K} \sum_{1 \leq i \neq j \leq K} \mathcal{D}(\mathcal{M}_i, \mathcal{M}_j) \\ & \text{s.t. } \forall i, \Delta_1 < |\mathcal{M}_i| < \Delta_2; \sum_i |\mathcal{M}_i| = n, \end{aligned} \quad (3)$$

where $\mathcal{D}(\cdot, \cdot)$ can be any distribution-related distance metric, n is the number of training samples in a single regime, Δ_1 , Δ_2 and K_0 are hyperparameters to avoid trivial partitions and over-splitting, \mathcal{M} denotes the subset of representations that corresponds to contiguous samples. Specifically, we choose the Deep Correlation Alignment (CORAL) Sun & Saenko (2016) to measure the temporal distribution similarity, *i.e.*, $\mathcal{D}(\cdot, \cdot) = \frac{1}{4q^2} \|C(\cdot) - C(\cdot)\|_F^2$, where q is the number of dimensions for each hidden state, $C(\cdot)$ denotes the second-order statistics (covariance matrix). The detailed optimization procedure is explained in Appendix A.

After the most diverse modes are obtained, the distribution of a regime can be efficiently preserved by selecting a small number of the most representative samples of each mode. The selection algorithm is shown in Algorithm 1, where we select samples that minimize CORAL to ensure that the selected small number of samples are well aligned/matched to each mode of MTS. By iterating all modes, we update the memory buffer as the union of all selected sample sets.

3.4 SEQUENTIAL TRAINING AND TESTING WITH SKI-CL

We briefly introduce the pipeline of continual MTS forecasting, as shown in Figure 1. At the training phase of i -th regime, the training objectives $\mathcal{L}_{\text{regime-}i}$ contains the objective for the current training samples, denoted as $\mathcal{L}_{\text{current}}$ and that for the memory of previous $i-1$ regimes, denoted as $\mathcal{L}_{\text{memory}}$, where $\mathcal{L}_{\text{regime-}i} = \mathcal{L}_{\text{current}} + \alpha\mathcal{L}_{\text{memory}}$, with α being the weight of memory loss.

After training, the structural knowledge of the current regime is saved in the structural memory and the training samples are selected to enrich the MTS memory as aforementioned. At the testing phase, SKI-CL is able to maintain the forecasting performance on the queries of testing samples from all regimes up to the current one, and accordingly recover the learned dependency structures informed by the structural knowledge of each regime (as shown in Figure 1(right)). Unlike other graph structure learning methods for continual MTS forecasting, our method is able to dynamically infer faithful dependency structures for existing and current regimes without accessing the memory buffer, which is more practical in real world applications.

Algorithm 1 Representation-Matching Sample Selection

- 1: **Input:** Sample representation H , hyperparameters Δ_1 , Δ_2 , K_0 , memory budget for a single regime N_m , the number of training samples in a single regime n
 - 2: Split H into modes $\mathcal{M}_1, \dots, \mathcal{M}_K$ by optimizing (3)
 - 3: Initialize an empty memory buffer S for a regime
 - 4: **for** $k \leftarrow 1$ to K **do**
 - 5: $n_{\text{sample}} \leftarrow 0$; $n_{\text{select}} \leftarrow N_m \times \frac{|\mathcal{M}_k|}{n}$; $s \leftarrow \{\}$
 - 6: **while** $n_{\text{sample}} \leq n_{\text{select}}$ **do**
 - 7: $i_{\text{selected}} \leftarrow \min_i \mathcal{D}(H_{s \cup i}, H_{\mathcal{M}_k}) \quad \triangleright i \notin s$
 - 8: $s = s \cup i_{\text{selected}}$; $n_{\text{sample}} += 1$
 - 9: $S = S \cup s$
 - 10: **Output:** The memory buffer S
-

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets To evaluate the performance of SKI-CL on continual MTS forecasting, we conduct experiments on three public benchmark MTS datasets including the traffic (Traffic-CL), solar energy (Solar-CL), and human activity recognition (HAR-CL), as well as one synthetic dataset based on Non-repeating Random Walk Denton (2005) that is used in Liu et al. (2022b) under continual learning setting. The statistics of these datasets are summarized in Table 1. For Traffic-CL and Solar-CL, the structural knowledge is the spatial proximity of the sensor/station. For HAR-CL, the partial structural knowledge is drawn from the domain-specific motion dynamics. For synthetic data, structural knowledge is the feature similarity of different variables. The edges from the structural knowledge are binary for Traffic-CL and HAR-CL, and continuous for Solar-CL and Synthetic-CL.

Baselines We compare SKI-CL with a number of dependency-modeling-based forecasting methods and commonly used continual learning methods to resolve the catastrophic forgetting issue in sequential training. The forecasting methods include statistical model VAR Lütkepohl (2005), ARIMA Box et al. (2015), and deep learning mod-

Table 1: Summary of Datasets

Dataset	Traffic-CL	Solar-CL	HAR-CL	Synthetic-CL
# of nodes	22	50	9	10
# of all time steps	106,848	52,560	600,576	24,000
# of regimes	7	5	4	4
Regime	year	state	activity	adjacency
Structure avail.	Completed	Completed	Partial	Completed

Table 2: Experiment Results for 12 Horizon Prediction. (Lower MAE and RMSE for AP mean better; When AP is comparable, lower MAE and RMSE for AF mean better.)

Model	Traffic-CL				Solar-CL				HAR-CL ($\times 10^{-2}$)				Synthetic-CL ($\times 10^{-2}$)			
	AP \downarrow		AF		AP \downarrow		AF		AP \downarrow		AF		AP \downarrow		AF	
	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
VAR _{seq}	88.19	126.01	58.38	80.58	167.30	534.42	205.27	658.80	19.59	28.38	1.93	2.19	22.34	32.70	9.18	13.54
ARIMA _{seq}	141.75	159.89	77.61	77.40	14.97	18.92	4.75	2.92	40.68	52.87	2.35	2.38	42.24	43.51	13.39	12.98
TCN _{seq}	16.88	28.67	3.77	6.83	2.03	4.84	0.06	0.24	14.85	23.42	3.60	5.06	4.30	4.90	0.66	0.99
TCN _{mir}	15.70	26.53	1.70	3.22	1.99	4.79	0.10	0.19	13.91	22.15	2.64	2.93	3.79	4.63	0.46	0.73
TCN _{herd}	15.55	26.21	1.49	2.81	2.01	4.82	0.13	0.23	13.87	22.08	2.05	2.88	3.72	4.61	0.35	0.67
TCN _{er}	15.51	26.23	1.46	2.80	1.98	4.73	-0.05	0.02	13.66	21.78	1.82	2.59	3.29	4.30	0.34	0.61
TCN _{der++}	15.46	25.68	1.33	2.49	1.95	4.69	-0.07	-0.02	13.56	21.55	1.69	2.28	3.00	4.00	0.28	0.32
ESG _{seq}	18.77	30.02	6.46	10.07	2.80	5.77	1.28	1.74	17.63	26.84	7.28	9.41	8.98	14.19	1.32	1.98
ESG _{mir}	18.24	29.83	5.02	8.25	2.03	4.83	0.25	0.49	17.25	26.63	4.01	5.21	8.95	13.91	1.21	1.81
ESG _{herd}	17.49	28.64	4.82	7.45	1.92	4.72	0.13	0.53	17.22	26.59	3.99	5.13	8.94	13.88	1.11	1.74
ESG _{er}	16.40	27.50	3.05	5.34	2.01	4.82	0.24	0.44	17.15	25.84	4.63	5.33	8.84	13.86	1.21	1.62
ESG _{der++}	17.40	29.21	4.01	6.97	1.91	4.57	0.09	0.21	16.20	24.32	5.18	6.00	8.81	13.77	1.02	1.42
GTS _{seq}	17.26	29.11	2.33	3.48	2.19	5.20	0.27	0.59	16.44	25.41	3.68	5.10	6.51	8.89	1.88	3.39
GTS _{mir}	17.17	29.08	2.13	3.31	2.15	5.16	0.14	0.68	15.83	24.85	3.27	4.91	6.44	8.57	1.31	2.86
GTS _{herd}	17.00	29.01	2.17	2.98	2.11	5.06	0.13	0.30	15.65	24.33	1.99	2.86	6.34	8.23	1.18	1.81
GTS _{er}	15.83	26.20	1.12	2.52	2.01	4.75	0.12	0.05	15.06	23.52	2.00	2.73	5.59	6.32	0.44	0.69
GTS _{der++}	15.84	26.05	1.15	2.33	1.94	4.57	-0.25	-0.19	14.80	23.01	1.52	1.88	5.43	6.67	0.23	0.30
MTGNN _{seq}	19.88	32.94	7.83	12.68	2.12	4.75	0.38	0.44	14.86	22.58	2.59	3.61	10.26	14.92	1.16	1.81
MTGNN _{mir}	18.01	31.84	5.03	8.97	2.00	4.73	0.21	0.40	14.59	22.52	2.24	3.53	8.92	12.91	1.07	1.33
MTGNN _{herd}	17.93	30.70	4.90	8.40	1.89	4.68	0.13	0.35	14.09	22.50	1.13	1.62	8.11	12.88	1.03	1.27
MTGNN _{er}	15.79	26.52	2.76	4.87	1.94	4.62	0.14	0.25	13.59	21.85	1.91	2.79	8.70	13.69	0.61	1.21
MTGNN _{der++}	15.40	25.99	2.22	4.10	1.90	4.57	0.06	0.14	13.57	21.75	1.63	2.40	8.63	13.51	0.50	0.92
Autoformer _{seq}	23.92	40.40	2.58	4.42	6.04	11.41	0.53	1.56	19.97	28.76	2.81	3.08	5.15	6.87	0.53	0.46
Autoformer _{mir}	23.85	40.13	2.21	4.07	5.95	11.18	0.87	1.14	18.57	27.93	1.45	2.88	5.12	6.73	0.32	0.37
Autoformer _{herd}	23.90	40.21	2.43	4.29	5.91	11.12	0.90	1.20	18.78	28.03	1.42	2.66	5.00	6.65	0.23	0.24
Autoformer _{er}	23.26	38.98	1.28	2.00	5.83	10.89	0.49	1.32	18.42	27.05	1.29	2.01	5.02	6.66	0.24	0.25
Autoformer _{der++}	22.11	37.34	1.08	1.92	5.34	9.29	0.21	1.14	18.12	26.91	1.12	1.85	4.97	6.61	0.18	0.23
PatchTST _{seq}	19.11	32.50	2.34	2.97	2.64	5.32	0.72	0.43	17.91	27.13	7.18	6.88	4.85	5.93	1.59	1.78
PatchTST _{mir}	19.04	32.23	2.28	2.79	2.61	5.30	0.70	0.40	17.82	26.89	6.82	4.81	4.83	5.86	1.55	1.72
PatchTST _{herd}	18.96	32.10	2.21	2.67	2.60	5.30	0.68	0.35	17.73	26.84	6.62	4.79	4.80	5.79	1.43	1.68
PatchTST _{er}	18.77	31.50	1.98	2.01	2.57	5.27	0.47	0.30	17.57	26.40	6.02	4.69	4.72	5.26	1.03	1.54
PatchTST _{der++}	18.53	31.34	1.75	1.98	2.53	5.17	0.43	0.28	17.12	26.13	5.79	4.32	4.64	5.13	0.83	0.88
DLinear _{seq}	19.69	32.75	2.91	2.83	3.47	6.56	1.17	1.12	17.32	26.31	2.71	3.43	4.81	5.81	1.64	1.57
DLinear _{mir}	19.37	32.25	2.17	2.59	3.45	6.51	1.02	1.01	16.87	26.12	2.67	3.01	4.79	5.73	1.47	1.40
DLinear _{herd}	19.53	32.40	2.25	2.68	3.41	6.50	1.03	1.00	16.83	25.81	2.57	2.91	4.77	5.70	1.59	1.46
DLinear _{er}	19.19	32.30	1.73	2.14	3.37	6.43	0.93	0.98	16.71	25.75	2.13	2.85	4.74	5.20	1.23	1.43
DLinear _{der++}	19.02	31.97	1.75	1.93	3.25	6.37	0.83	0.79	16.58	25.47	1.92	2.77	4.21	4.88	1.12	1.13
TimesNet _{seq}	17.77	29.91	3.13	6.93	3.92	7.18	1.46	2.51	18.38	27.61	4.33	5.15	5.18	6.13	1.72	2.03
TimesNet _{mir}	17.53	29.61	2.44	5.32	3.77	7.15	1.22	1.57	18.27	27.59	4.01	5.08	5.12	6.05	1.69	1.97
TimesNet _{herd}	17.38	29.53	2.56	5.83	3.83	7.10	1.03	1.44	18.01	27.53	3.46	5.03	5.10	6.03	1.68	1.95
TimesNet _{er}	17.25	29.33	1.97	4.19	3.55	7.02	0.42	0.91	17.84	27.07	3.28	4.01	4.93	5.90	1.42	1.90
TimesNet _{der++}	17.13	29.28	1.56	4.02	3.45	6.55	0.37	0.90	17.73	26.86	3.11	3.87	4.81	5.88	1.32	1.78
SKI-CL _{seq}	17.30	29.38	4.38	7.80	2.02	4.73	0.30	0.50	14.73	23.31	3.91	5.07	4.70	5.85	1.97	3.46
SKI-CL _{mir}	15.77	26.32	1.95	3.47	1.98	4.69	0.35	0.56	13.65	21.77	2.61	3.82	4.53	5.01	1.67	2.21
SKI-CL _{herd}	15.45	25.73	1.82	3.28	2.00	4.70	0.33	0.52	13.71	21.82	2.53	3.67	4.44	4.82	1.23	2.02
SKI-CL _{er}	15.43	25.60	1.69	2.92	1.95	4.67	0.11	0.23	13.58	21.57	1.82	2.50	3.39	4.52	0.30	0.38
SKI-CL _{der++}	15.39	25.57	1.63	2.87	1.91	4.60	0.10	0.21	13.50	21.47	1.74	2.42	3.33	4.43	0.28	0.30
SKI-CL	15.24	25.33	1.50	2.71	1.75	4.46	0.09	0.06	13.41	21.30	1.64	2.08	3.24	4.24	0.15	0.23

els including TCN Bai et al. (2018), LSTNet Lai et al. (2018), STGCN Yu et al. (2018), MTGNN Wu et al. (2020), AGCRN Bai et al. (2020), GTS Shang & Chen (2021) and ESG Ye et al. (2022), StemGNN Cao et al. (2020), Autoformer Wu et al. (2021), PatchTST Nie et al. (2022), Dlinear Zeng et al. (2023), TimesNet Wu et al. (2022) where STGCN and GTS use structural knowledge and ESG learns a dynamic graph in MTS modeling. All forecasting methods are first evaluated on sequential training without any countermeasures (denoted as seq). The continual learning methods employ the memory-replay-based methods including the herding method (denoted as herd) Rebuffi et al. (2016), the randomly replay training samples (denoted as er), a DER++ method Buzzega et al. (2020) that enforces a L_2 knowledge distillation loss on the previous logits (denoted as der++) and a MIR method that selects samples with highest forgetting for experience replay (denoted as mir). For the proposed SKI-CL, we also evaluate our proposed representation-matching sample selection scheme.

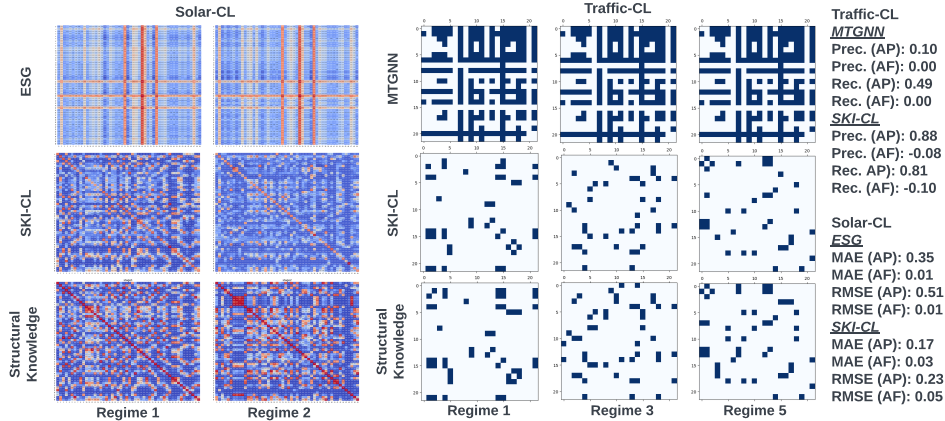


Figure 3: The Structure Visualizations on Traffic-CL and Solar-CL datasets.

Evaluation Metrics We adopt two metrics based on Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) to evaluate the performance on continual MTS forecasting, *i.e.*, the Average Performance (AP) and Average Forgetting (AF) Lopez-Paz & Ranzato (2017); Zhang et al. (2022a). The AP at i -th regime is defined as $AP = \sum_{j=1}^i P_{i,j}/i$ for $i \geq 1$, where $P_{i,j}$ denotes the performance on regime j after the model has been sequentially trained from stage 1 to i . Similarly, the Average Forgetting is defined as $AF = \sum_{j=1}^{i-1} (P_{i,j} - P_{j,j})/(i-1)$ for $i \geq 2$. Besides the forecasting performance, we also evaluate AP and AF on the learned dependency structures, where the average precision (Prec.) and average recall (Rec.) are used for a binary graph, MAE and RMSE are used for a continuous graph. More details of datasets, baselines, and evaluations are provided in Appendix C.

4.2 PERFORMANCE EVALUATION FOR CONTINUAL MTS FORECASTING

In this paper, we focus on a multi-horizon continual forecasting task. Table 2 summarizes the experiment results of [selected baselines \(full experiment results are shown in Appendix Table 7\)](#) and our proposed SKI-CL method with variants for 12 horizon predictions (the corresponding standard deviations and results for different horizons are provided in Table 9 and 8, Appendix F respectively.). Based on the results, we show that statistical methods, such as VAR and ARIMA, cannot perform well under continual learning settings and all other forecasting methods suffer from obvious performance degradation (AF) in sequential training, indicating the existence of catastrophic forgetting phenomena when regime shifts. In addition, the memory-replay-based methods generally alleviate the forgetting issues with better APs (lower RMSE and MAE) and smaller relative AFs. SKI-CL and its variants (SKI-CL_{er} and SKI-CL_{der++}) consistently achieve the best or the second-best APs, showing advantages over other models equipped with memory-replay-based methods (*e.g.*, MTGNN_{der++}, TCN_{er}, GTS_{der++}) since SKI-CL explicit imposes a regularization to enforce the consistency between the learned dynamic graphs and the structure knowledge. These observations demonstrate that learning a dynamic structure is beneficial for nonstationary MTS modeling, and the structural knowledge helps to characterize the general variable behaviors in each regime. Even partially observed structural knowledge can serve as a valid reference to learn the dependency structures. Finally, the proposed SKI-CL achieves better APs or less AFs when the AP metrics are comparable, suggesting the superiority of our representation-matching scheme that maximizes the coverage of dependencies and temporal dynamics.

4.3 PRESERVING FAITHFUL DEPENDENCY STRUCTURES

We also evaluate how the baselines and our method preserve the learned structures that are highly correlated to the structural knowledge. Specifically, we compare SKI-CL with GTS on Traffic-CL dataset, and ESG on Solar-CL dataset to investigate the binary edges and continuous edges, respectively. The results of the learned structures and structural knowledge are shown in Figure 3, where the average performance and average forgetting are also annotated. It is clear that SKI-CL is able to alleviate the forgetting of the learned structures on both datasets at the testing phase, while the baselines fail to model the nonstationarity of MTS in different regimes. This observation also

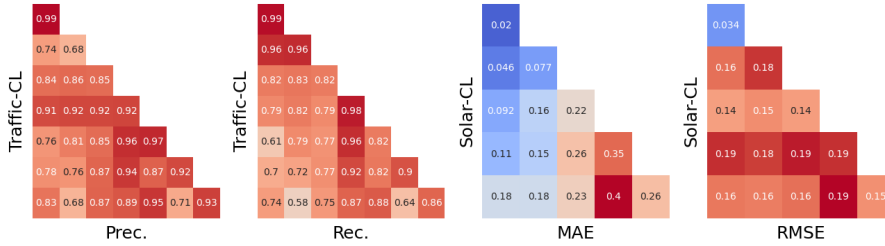


Figure 4: The Structure Similarity Matrices on Traffic-CL and Solar-CL datasets.

Table 3: Analysis of λ on Traffic-CL

λ	Forecasting Performance				Structure Similarity			
	AP		AF		AP		AF	
	MAE	RMSE	MAE	RMSE	Prec.	Recall	Prec.	Recall
0.01	15.42	26.08	2.30	4.13	0.52	0.46	-0.52	-0.57
0.1	15.39	25.97	2.12	4.07	0.84	0.75	-0.11	-0.19
0.5	15.33	25.68	1.68	3.07	0.85	0.78	-0.08	-0.11
1	15.22	25.31	1.53	2.74	0.84	0.80	-0.10	-0.12
2	15.27	25.51	1.70	3.10	0.85	0.79	-0.08	-0.14
5	15.31	25.63	2.04	3.56	0.83	0.72	-0.16	-0.24

Table 4: Analysis of Memory Budget on Traffic-CL

Ratio	Forecasting Performance				Structure Similarity			
	AP		AF		AP		AF	
	MAE	RMSE	MAE	RMSE	Prec.	Recall	Prec.	Recall
0.01	15.22	25.31	1.53	2.74	0.84	0.80	-0.10	-0.12
0.05	14.49	24.35	1.16	2.03	0.86	0.79	-0.10	-0.13
0.1	14.44	24.03	0.77	1.23	0.84	0.78	-0.07	-0.13
0.2	14.25	23.74	0.85	1.34	0.89	0.80	-0.06	-0.14
0.5	14.18	23.06	0.55	0.79	0.90	0.81	-0.05	-0.09

suggests the importance of dynamic structure learning and the incorporation of structural knowledge, to maintain a faithful structure at each regime. We further visualize the similarity matrices between the structures inferred by SKI-CL and structural knowledge, as shown in Figure 4. Each cell corresponds to the aforementioned $P_{i,j}$, *i.e.*, the performance on regime j after the model has been sequentially trained from stage 1 to i , where i and j denote the row number and column number, respectively. We observe that the inferred structures at the testing phases of each regime still reveal similarities to the structural knowledge, by comparing the values of each row with the diagonal ones.

We emphasize that we don’t intend to use structural knowledge as a ground truth. We have demonstrated that exploiting structural knowledge helps to reduce the performance degradation between consecutive regimes. Besides, it is beneficial to have the model aligned with the structural knowledge for a better interpretation of each regime. We provide a case study on the Synthetic-CL dataset as a demonstration in Figure 9, Appendix G . More visualization results are provided in Appendix D.

4.4 HYPERPARAMETER ANALYSIS

We perform experiments on the Traffic-CL dataset to validate the effectiveness and sensitivity of two key hyperparameters in SKI-CL, the weight of structure regularizer λ (1 by default) and the memory budget (sampling ratio) at each regime (0.01 by default). As shown in Table 3, within a small range, the model is relatively stable in terms of λ , resulting in similar average performance and average forgetting on forecasting performance as well as the inferred structure similarity. We also study the model performance when the memory budget varies in Table 4. We can observe that a larger memory budget can achieve better forecasting performance and less forgetting. Meanwhile, the structure similarity is relatively stable regarding the memory budget. The exploration of other hyperparameter settings is provided in Appendix E.

5 CONCLUSION

In this paper, we propose a novel Structural Knowledge Informed Continual Learning (SKI-CL) framework to perform continual MTS forecasting, which exploits the structural knowledge to characterize the dynamic variable dependencies within each regime. We first design a dynamic graph learning model for fine-grained dynamic MTS modeling and impose a regularization scheme that ensures the consistency between the learned structure and the variable dependencies revealed by structural knowledge. We further alleviate the catastrophic forgetting by developing a representation-matching memory replay scheme, which maximizes the temporal coverage of MTS data to efficiently preserve each regime’s underlying temporal dynamics and dependency structure. Experiments on one synthetic dataset and three real-world benchmark datasets demonstrate the effectiveness and advantages of the proposed SKI-CL on continual MTS forecasting tasks.

REFERENCES

- Davide Anguita, Alessandro Ghio, Luca Oneto, Francesc Xavier Llanas Parra, and Jorge Luis Reyes Ortiz. Energy efficient smartphone-based activity recognition using fixed-point arithmetic. *Journal of universal computer science*, 19(9):1295–1314, 2013a.
- Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, Jorge Luis Reyes-Ortiz, et al. A public domain dataset for human activity recognition using smartphones. In *Esann*, volume 3, pp. 3, 2013b.
- Lei Bai, Lina Yao, Can Li, Xianzhi Wang, and Can Wang. Adaptive graph convolutional recurrent network for traffic forecasting. *Advances in Neural Information Processing Systems*, 33:17804–17815, 2020.
- Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv:1803.01271*, 2018.
- Anastasia Borovykh, Sander Bohte, and Cornelis W Oosterlee. Conditional time series forecasting with convolutional neural networks. *arXiv preprint arXiv:1703.04691*, 2017.
- George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. *Time series analysis: forecasting and control*. John Wiley & Sons, 2015.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930, 2020.
- Defu Cao, Yujing Wang, Juanyong Duan, Ce Zhang, Xia Zhu, Congrui Huang, Yunhai Tong, Bixiong Xu, Jing Bai, Jie Tong, et al. Spectral temporal graph neural network for multivariate time-series forecasting. *Advances in Neural Information Processing Systems*, 33:17766–17778, 2020.
- Chao Chen, Karl Petty, Alexander Skabardonis, Pravin Varaiya, and Zhanfeng Jia. Freeway performance measurement system: mining loop detector data. *Transportation Research Record*, 1748(1): 96–102, 2001.
- Weijun Chen, Yanze Wang, Chengshuo Du, Zhenglong Jia, Feng Liu, and Ran Chen. Balanced graph structure learning for multivariate time series forecasting, 2022.
- Xu Chen, Junshan Wang, and Kunqing Xie. Trafficstream: A streaming traffic flow forecasting framework based on graph neural networks and continual learning.
- Anne Denton. Kernel-density-based clustering of time series subsequences using a continuous random-walk noise model. In *Fifth IEEE International Conference on Data Mining (ICDM'05)*, pp. 8–pp. IEEE, 2005.
- Yuntao Du, Jindong Wang, Wenjie Feng, Sinno Pan, Tao Qin, Renjun Xu, and Chongjun Wang. Adarnn: Adaptive learning and forecasting of time series. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pp. 402–411, 2021.
- Ziheng Duan, Haoyan Xu, Yida Huang, Jie Feng, and Yueyang Wang. Multivariate time series forecasting with transfer entropy graph. *Tsinghua Science and Technology*, 28(1):141–149, 2022.
- Aurora Gonzalez-Vidal, Fernando Jimenez, and Antonio F Gomez-Skarmeta. A methodology for energy multivariate time series forecasting in smart buildings based on feature selection. *Energy and Buildings*, 196:71–82, 2019.
- Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.
- Silviu Guiasu and Abe Shenitzer. The principle of maximum entropy. *The mathematical intelligencer*, 7:42–48, 1985.
- Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, and Huaiyu Wan. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 922–929, 2019.

- Vibhor Gupta, Jyoti Narwariya, Pankaj Malhotra, Lovekesh Vig, and Gautam Shroff. Continual learning for multivariate time series tasks with variable input dimensions. In *2021 IEEE International Conference on Data Mining (ICDM)*, pp. 161–170. IEEE, 2021.
- Yujiang He and Bernhard Sick. Clear: An adaptive continual learning framework for regression tasks. *AI Perspectives*, 3(1):1–16, 2021.
- Elvin Isufi, Andreas Loukas, Nathanael Perraudin, and Geert Leus. Forecasting time series with varma recursions on graphs. *IEEE Transactions on Signal Processing*, 67(18):4870–4885, 2019.
- Bo Jin, Haoyu Yang, Leilei Sun, Chuanren Liu, Yue Qu, and Jianing Tong. A treatment engine by predicting next-period prescriptions. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1608–1616, 2018.
- Zixuan Ke and Bing Liu. Continual learning of natural language processing tasks: A survey. *arXiv preprint arXiv:2211.12701*, 2022.
- Mahdi Khodayar and Jianhui Wang. Spatio-temporal graph deep neural network for short-term wind speed forecasting. *IEEE Transactions on Sustainable Energy*, 10(2):670–681, 2018.
- Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*, 2021.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pp. 95–104, 2018.
- Jaehoon Lee, Chan Kim, Gyumin Lee, Haksoo Lim, Jeongwhan Choi, Kookjin Lee, Dongeun Lee, Sanghyun Hong, and Noseong Park. Time series forecasting with hypernetworks generating parameters in advance. *arXiv preprint arXiv:2211.12034*, 2022.
- Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. *arXiv preprint arXiv:1707.01926*, 2017.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- Haozhe Lin, Yushun Fan, Jia Zhang, and Bing Bai. Rest: Reciprocal framework for spatiotemporal-coupled predictions. In *Proceedings of the Web Conference 2021*, pp. 3136–3145, 2021.
- Jiahao Liu, Dongsheng Li, Hansu Gu, Tun Lu, Peng Zhang, and Ning Gu. Parameter-free dynamic graph embedding for link prediction. *Advances in Neural Information Processing Systems*, 35: 27623–27635, 2022a.
- Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *International conference on learning representations*, 2021.
- Yijing Liu, Qinxian Liu, Jian-Wei Zhang, Haozhe Feng, Zhongwei Wang, Zihan Zhou, and Wei Chen. Multivariate time-series forecasting with temporal polynomial graph neural networks. In *Advances in Neural Information Processing Systems*, 2022b.
- David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. In *NIPS*, 2017.
- Yonghong Luo, Chao Lu, Lipeng Zhu, and Jie Song. Data-driven short-term voltage stability assessment based on spatial-temporal graph convolutional network. *International Journal of Electrical Power & Energy Systems*, 130:106753, 2021.

- Helmut Lütkepohl. *New introduction to multiple time series analysis*. Springer Science & Business Media, 2005.
- Rohan Money, Joshin Krishnan, and Baltasar Beferull-Lozano. Online non-linear topology identification from graph-connected time series. In *2021 IEEE Data Science and Learning Workshop (DSLW)*, pp. 1–6. IEEE, 2021.
- Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 4602–4609, 2019.
- Alberto Natali, Elvin Isufi, Mario Coutino, and Geert Leus. Learning time-varying graphs from online data. *IEEE Open Journal of Signal Processing*, 3:212–228, 2022.
- William T Ng, K Siu, Albert C Cheung, and Michael K Ng. Expressing multivariate time series as graphs with time series attention transformer. *arXiv preprint arXiv:2208.09300*, 2022.
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*, 2022.
- Namyong Park, Ryan Rossi, Eunye Koh, Iftikhar Ahamath Burhanuddin, Sungchul Kim, Fan Du, Nesreen Ahmed, and Christos Faloutsos. Cgc: Contrastive graph clustering for community detection and tracking. In *Proceedings of the ACM Web Conference 2022*, pp. 1115–1126, 2022.
- Quang Pham, Chenghao Liu, Doyen Sahoo, and Steven CH Hoi. Learning fast and slow for online time series forecasting. *arXiv preprint arXiv:2202.11672*, 2022.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, G. Sperl, and Christoph H. Lampert. icarl: Incremental classifier and representation learning. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5533–5542, 2016.
- David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks.
- Chao Shang and Jie Chen. Discrete graph structure learning for forecasting multiple time series. In *Proceedings of International Conference on Learning Representations*, 2021.
- Shun-Yao Shih, Fan-Keng Sun, and Hung-yi Lee. Temporal pattern attention for multivariate time series forecasting. *Machine Learning*, 108(8):1421–1441, 2019.
- Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *ECCV Workshops*, 2016.
- Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alexander Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. 2016.
- Jindong Wang, Yiqiang Chen, Shuji Hao, Xiaohui Peng, and Lisha Hu. Deep learning for sensor-based activity recognition: A survey. *Pattern recognition letters*, 119:3–11, 2019.
- Yiwei Wang, Yujun Cai, Yuxuan Liang, Henghui Ding, Changhu Wang, Siddharth Bhatia, and Bryan Hooi. Adaptive data augmentation on temporal graphs. *Advances in Neural Information Processing Systems*, 34:1440–1452, 2021.
- Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 139–149, June 2022.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems*, 34:22419–22430, 2021.

- Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. *arXiv preprint arXiv:2210.02186*, 2022.
- Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. Connecting the dots: Multivariate time series forecasting with graph neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 753–763, 2020.
- Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Junchen Ye, Zihan Liu, Bowen Du, Leilei Sun, Weimiao Li, Yanjie Fu, and Hui Xiong. Learning the evolutionary and multi-scale graph structure for multivariate time series forecasting. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '22*, pp. 2296–2306, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393850. doi: 10.1145/3534678.3539274. URL <https://doi.org/10.1145/3534678.3539274>.
- Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: a deep learning framework for traffic forecasting. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pp. 3634–3640, 2018.
- Bakht Zaman, Luis Miguel Lopez Ramos, Daniel Romero, and Baltasar Beferull-Lozano. Online topology identification from vector autoregressive time series. *IEEE Transactions on Signal Processing*, 69:210–225, 2020.
- Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 11121–11128, 2023.
- Liheng Zhang, Charu Aggarwal, and Guo-Jun Qi. Stock price prediction via discovering multi-frequency trading patterns. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 2141–2149, 2017.
- Xikun Zhang, Dongjin Song, and Dacheng Tao. Cglb: Benchmark tasks for continual graph learning. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022a.
- Xikun Zhang, Dongjin Song, and Dacheng Tao. Hierarchical prototype networks for continual graph representation learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022b.
- YiFan Zhang, Qingsong Wen, Xue Wang, Weiqi Chen, Liang Sun, Zhang Zhang, Liang Wang, Rong Jin, and Tieniu Tan. Onenet: Enhancing time series forecasting models under concept drift by online ensembling. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of AAAI*, 2021.
- Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International Conference on Machine Learning*, pp. 27268–27286. PMLR, 2022.

A REPRESENTATION-BASED DISTRIBUTION CHARACTERIZATION

A.1 OPTIMIZING THE OBJECTIVE FOR COVERAGE MAXIMIZATION

Inspired by the principle of maximum entropy, we perform a distribution characterization by splitting the MTS data into the most diverse modes on the representation space, which incorporates variable dependence and temporal information. The distribution characterization splits the hidden representation by solving a constrained optimization problem whose objective is formulated as:

$$\begin{aligned} & \max_{0 < K \leq K_0} \max_{n_1, \dots, n_K} \frac{1}{K} \sum_{1 \leq i \neq j \leq K} \mathcal{D}(\mathcal{M}_i, \mathcal{M}_j) \\ & \text{s.t. } \forall i, \Delta_1 < |\mathcal{M}_i| < \Delta_2; \sum_i |\mathcal{M}_i| = n, \end{aligned} \quad (4)$$

where $\mathcal{D}(\cdot, \cdot)$ can be any distribution-related distance metric, n is the number of training samples in a single regime, Δ_1 , Δ_2 and K_0 are hyperparameters to avoid trivial partitions and over-splitting, \mathcal{M} denotes the subset of representations that correspond to contiguous samples. The optimization problem can be computationally intractable and the closed-form solution may not exist. We adopt the greedy algorithm proposed by Du et al. (2021). First, we obtain the ordered representation set from the trained model for single regime data. Then, for efficient computation, we evenly split the representation into N parts and randomly search the value of K in $\{2, 3, 4 \dots N - 1\}$. Denote the start and the end index of the representation by A and B respectively. We first consider $K = 2$ by choosing the first splitting point C from all candidate splitting points via maximizing the distance metric $\mathcal{D}(\mathcal{M}_{AC}, \mathcal{M}_{CB})$. After C is determined, we then consider $K = 3$ and use the same strategy to select another point D. A similar strategy is applied until the number of representation modes is obtained.

B RELATION TO OTHER RESEARCH TOPICS

Our method characterizes dynamic variable dependencies of MTS data within each regime to perform continual MTS forecasting tasks. The dynamic variable dependencies modeling and continual adaptation nature relate our method to two research topics, *i.e.*, the dynamic graph learning and learning with temporal drift of MTS. Dynamic graph learning focuses on adapting to the changes in the explicit graph structure and possibly features over time, and leverage it to facilitate downstream tasks (*e.g.*, node classification Liu et al. (2022a), link prediction Wang et al. (2021), community detection Park et al. (2022)). In contrast, our proposed continual multivariate time series forecasting method needs to discover underlying dependencies structure of variables over different regimes (stages), by capturing the temporal dynamics of MTS data. On the other hand, temporal drift of MTS represents the phenomenon where the data distribution of the target MTS changes over time in unforeseen ways, where the approaches are often more reactive and focused on adapting to changes in data distribution for a specific task Du et al. (2021); Kim et al. (2021); Lee et al. (2022). Based on the notion, it is important to emphasize the continual adaptation nature of continual learning setting regarding the temporal drifts, where the catastrophic forgetting happens due to the shifts over multiple regimes. That being said, continual MTS forecasters need to handle a variety of regimes, not just adapt to changes in data distribution for a single one, which is to some degree more realistic and challenging. Temporal drift represents a specific challenge within this broader spectrum of adaptation.

Moreover, while the aforementioned topics as well as our study deal with the evolving data, the learning objectives can be very different. For the temporal drift methods, the objective is to adapt the forecaster to these new MTS patterns by effectively detecting, responding to, and learning from these changes. Similarly, the focus of dynamic graph learning is to adapt its model to cope with the structural and possible feature changes. However, the main focus this paper is to prevent the model from forgetting previously learned knowledge (coupled temporal dynamics and variable dependencies in our scope) when adapting to new MTS distributions.

C DATASETS, BASELINES AND EVALUATIONS FOR CONTINUAL MTS FORECASTING

C.1 DATASETS

Traffic-CL Following the fashion in PEMS3-Stream Chen et al., we construct the Traffic-CL dataset based on the PEMS3 benchmark for continual MTS forecasting tasks. The PEMS3 benchmark data was collected by the Performance Measurement System in California Chen et al. (2001) in real-time by every 30 seconds and further aggregated to 5-min granularity. The PEMS3-Stream dataset contains traffic data from 2011 to 2017. Specifically, data within a month period from July 10th to August 9th from every year was selected, where the traffic network keeps expanding from year to year. The adjacency matrix for τ -th year is extracted by applying a Gaussian kernel to the spatial all pairwise distances between two traffic sensors, as shown in Equation 5.

$$A_{\tau}[i, j] = \begin{cases} \exp\left(-\frac{d_{ij}^2}{\sigma_d^2}\right), & i \neq j \quad \text{and} \quad d_{ij} \leq \epsilon \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

where d_{ij} denotes the spatial distance between sensor i and j . σ_d and ϵ are the standard deviation and a predefined threshold (controlling the sparsity of the adjacency matrix, set as 1), respectively.

Based on the constructed PEMS3-Stream dataset, we make the following modifications to further simulate distinct regimes in the setting of continual forecasting. For each year, we rank and select the top 100 traffic sensors with the largest node degrees, based on which we randomly select 22 sensors as a set representing a part of the major traffic. Next, we transform the continuous adjacency weights to binary ones by a threshold, and use it as the structural knowledge. As such, each regime is represented by a different portion of a temporally expanding traffic network from different years, with a binarized structural prior and MTS data defined accordingly. The input horizon for the Traffic-CL dataset is 12.

Solar-CL We build our continual MTS forecasting dataset based on the database for NREL’s Solar Power Data for Integration Studies¹, which contains 5-minute solar power data for near 6,000 simulated photovoltaic power plants in the United States for the year 2006. Note that the data in Alabama with a 10-minute granularity is also known as the commonly used Solar dataset in many existing MTS studies Wu et al. (2020); Lai et al. (2018); Cao et al. (2020); Liu et al. (2022b).

We construct different regimes by states with different average annual sunlight levels (measured by kJ/m^2). Based on the statistics² and the aggregated MTS data at 10-minute, we select five states, Massachusetts/MA ($3944 \text{ kJ}/\text{m}^2$) - Arizona/AZ ($5755 \text{ kJ}/\text{m}^2$) - North Carolina/NC ($4456 \text{ kJ}/\text{m}^2$) - Texas/TX ($5137 \text{ kJ}/\text{m}^2$) - Washington/WA ($3467 \text{ kJ}/\text{m}^2$) as five regimes representing different sunlight patterns in different spatial locations. For each state, 50 photovoltaic power plants are randomly selected, where the spatial information is also used as a valid structural prior, as plants that are geographically close share similar weather and sunlight conditions at a local level. Specially, we first extract the longitude and latitude for each plant and calculate the pairwise geographic distances among all plant pairs. Based on all pairwise distances, we generate the adjacency matrix by applying a Gaussian kernel in Equation 3, where we set $\epsilon = \infty$, indicating a fully connected continuous graph. The input horizon for the Solar-CL dataset is 24.

HAR-CL We leverage the class boundaries in MTS classification data to construct different regimes in forecasting tasks. Specifically, we build our continual forecasting dataset based on a commonly used MTS Classification benchmark, the Human Activity Recognition (HAR) dataset Anguita et al. (2013b;a) in the UCI database³, where the data is collected from a group of 30 volunteers from 19-48 years, wearing a Samsung smartphone on the waist and performing six activities of daily living (Walking, Walking upstairs, Walking downstairs, Standing, Sitting, Lying). Each MTS sequence contains 128 time steps and 9 variables that are recorded based on the accelerometer and gyroscope, including the linear accelerations, the angular velocities, and total accelerations along the X-Y-Z axis. The detailed setup for data collection can be found in Figure 1 of Anguita et al. (2013a).

¹<https://www.nrel.gov/grid/solar-power-data.html>

²<https://worldpopulationreview.com/state-rankings/sunniest-states>

³<https://archive.ics.uci.edu/ml/datasets/Human+Activity+Recognition+Using+Smartphones>

We notice that different human activities naturally form distinct regimes with unique temporal dynamics, where deep learning methods easily achieve over 94% classification accuracy Wang et al. (2019). Motivated by this fact, we construct the continual HAR forecasting dataset considering the following details. Firstly, we select Walking, Walking upstairs, Walking downstairs, and Lying as four different regimes in our task. Secondly, each regime contains over one thousand sequences that are not always temporally connected due to different volunteers. As such, we iterate all sequences and construct the input-output MTS windows in each 128-step sequence, after which all windows are stacked as the training/validation/testing data in one regime. Secondly, we try to build the structural knowledge for each regime based on the domain knowledge. We’ve already known that linear accelerations are highly correlated to total accelerations along each axis, regardless of the activities, but the dependencies among other variables are not very clear. Therefore, we examine the Pearson correlations of all training sequences for each regime, where the mean of correlation matrices demonstrate distinct patterns, and validate the strong correlations between linear and total accelerations. However, the standard deviations are only small at the diagonal and the aforementioned entries, suggesting a varying and uncertain structure for other variable pairs. To this end, we check a small (15-th) percentile of each entry in the absolute correlation matrix, and apply a threshold to get a binary mask representing the variable dependencies that we are confident of. Note that we only regularize the learned structures to be consistent with the partially observed structural knowledge at the masked entries. As such, we simulate the structured scenario when we are not aware of the completed structural knowledge of MTS. The input for the HAR-CL dataset has 12 steps.

Synthetic-CL Lastly, we generate the synthetic data based on Non-repeating Random Walk Denton (2005), which is used in Liu et al. (2022b) for the evaluation of the learned graph structure in a single regime. Next, we introduce how to generate the MTS data for continual forecasting tasks.

Firstly, we describe how to generate the MTS data step by step. At time step t , given a dynamic weighted adjacency matrix $\mathbf{W}^{(t)}$, $\mathbf{X}^{(t)} = \mathcal{N}(\mathbf{W}^{(t-1)}\mathbf{X}^{(t-1)}, \sigma) \in \mathbb{R}^{N \times 1}$, where \mathcal{N} denotes the Gaussian distribution, $\sigma \in \mathbb{R}$ controls the variance, N denotes the number of variables, and $\mathbf{X}^{(0)}$ is randomly initialized from the set $[-1, -0.5, 0.5, 1]$. Secondly, we describe how the $\mathbf{W}^{(t)}$ is generated and how to construct different regimes based on $\mathbf{W}^{(\cdot)}$. Assuming there are L total time steps, we define $\mathbf{W}^{(t)}$ as one of S constant matrices $(\mathbf{G}^{(1)}, \dots, \mathbf{G}^{(S)})$, where $\mathbf{G}^{(\cdot)}$ is a Laplacian of sparsified random adjacency matrix with sparsity δ , spanning $\lfloor L/S \rfloor$ time steps ($\lfloor \cdot \rfloor$ denotes the floor function). As such, each regime is represented by MTS data with time steps from $(i-1) \times \lfloor L/S \rfloor$ to $i \times \lfloor L/S \rfloor$, with the corresponding weighted adjacency matrix $\mathbf{W}^{(t)} = \mathbf{G}^{(i)}$. Specifically, we set the number of variable $N = 10$, the total time steps $L = 24,000$, the number of regimes $S = 4$, the standard deviation of noise $\sigma = 0.01$, and the matrix sparsity $\delta = 0.1$.

There are two main differences between the synthetic setting of Liu et al. (2022b) and our work despite the setting of continual learning. Firstly, we don’t reinitialize the value of each variable when the dynamic weighted adjacency matrix transits to a different one in order to preserve the temporal continuity of MTS data across different regimes. Secondly, the evaluation of graph structure learning is also different due to the forecasting setting. In this particular synthetic dataset, the dynamic weighted adjacency matrix $\mathbf{W}^{(t)}$ describes the data generation process at the single-step level, which can be treated as the ground truth if the non-linear part of $\mathbf{W}^{(t)}$ in the model learns an identity mapping with Gaussian noise. As the graph learned in Liu et al. (2022b) is under the setting of single-step prediction, the $\mathbf{W}^{(t)}$ itself is a reasonable reference for evaluation. In our cases of performing multi-horizon forecasting, the matrix $\mathbf{W}^{(t)}$ raised to a higher power can also demonstrate how the dynamic is propagated in a sequence. In our exploration, we don’t assume $\mathbf{W}^{(t)}$ is explicitly given as the structural prior. Instead, we exploit the Pearson correlation of the generated MTS data and formulate a binary structural prior based on strong absolute correlations, where we will examine if the learned graph structure is able to reveal the variable interactions in $\mathbf{W}^{(t)}$. The input horizon for the Synthetic-CL dataset is 12.

C.2 BASELINES

In this part, we introduce the state-of-the-art baseline methods evaluated in our paper.

- TCN Bai et al. (2018): Temporal convolution networks (TCN) models the temporal causality using causal convolution and do not involve structural dependence modeling.

- LSTNet Lai et al. (2018): LSTNet leverages the Convolution Neural Network (CNN) with a kernel spanning the variable dimension to extract short-term local variable dependencies, and the Recurrent Neural Network (RNN) to discover long-term patterns based on the extracted dependency patterns for MTS forecasting.
- STGCN Yu et al. (2018): STGCN jointly captures the spatial-temporal patterns by stacking spatial graph convolution layers that perform graph convolution using continuous structural knowledge and temporal-gated convolution layers that capture temporal dynamics based on the yielded spatial representations.
- MTGNN Wu et al. (2020): MTGNN learns a parameterized graph with top-k connections for each node, and performs mix-hop propagation for graph convolution and dilated inception for temporal convolution. The parameterized graph is purely optimized based on the forecasting objective. At the testing stage, the inferred graph is static due to the fixed parameters.
- AGCRN Bai et al. (2020): AGCRN models the dependencies graph structure as a product of trainable node embedding and performs graph convolution in the recurrent convolution layer for MTS forecasting. The node embedding and yielded graph are purely optimized based on the forecasting objective. At the testing stage, the inferred graph is static due to the fixed parameters.
- GTS Shang & Chen (2021): GTS infers steady node representations and global node relations from entire training MTS data. The learned dependency structure is used in Diffusion Convolution Recurrent Neural Networks (DCRNN) for MTS forecasting. The parameterized graph is optimized based on the forecasting objective as well as the regularization based on binary structure priors. At the testing stage, the graph is sampled from learned binary edge distributions.
- ESG Ye et al. (2022): ESG learns evolving and scale-specific node relations from features extracted from MTS data. A series of dynamic graphs representing dynamic correlations are utilized in sequential graph convolution and temporal convolution. The dynamic graphs are learned via the optimization of feature extraction layers based on the forecasting objective. At the testing stage, the graphs are dynamics inferred based on each MTS input window.
- StemGNN Cao et al. (2020): The Spectral Temporal Graph Neural Network (StemGNN) is a Graph-based multivariate time-series forecasting model, which jointly learns temporal dependencies and inter-series correlations in the spectral domain, by combining Graph Fourier Transform (GFT) and Discrete Fourier Transform (DFT).
- Autoformer Wu et al. (2021): Autoformer is a Transformer-based model using decomposition architecture with an Auto-Correlation mechanism to capture cross-time dependency for forecasting.
- PatchTST Nie et al. (2022): PatchTST model uses channel-independent and patch techniques to tokenize input time series and perform time series forecasting by utilizing the vanilla Transformer encoders.
- DLinear Zeng et al. (2023): DLinear adopts trend-seasonal components decomposition techniques for time series data and applies MLP-based architectures for time series forecasting.
- TimesNet Wu et al. (2022): TimesNet model leverages intricate temporal patterns by exploring the multi-periodicity of time series and capture the temporal 2D-variations in 2D space using transformer based backbones.

The differences between the graph learning baselines and our method are illustrated in Table 5. Compared with the state-of-the-art method, our proposed SKI-CL backbone model learns a dynamic graph for nonstationary MTS modeling, which is also capable of incorporating structural knowledge with different forms and availability scenarios so as to characterize the dependency structure and temporal dynamics of each regime.

Training Details The dynamic graph inference module consists of 3 stacked 2D convolutional layers. Using C_{in} , C_{out} to denote the number of channels coming in and out, the parameters of these convolutional layers are [$C_{in} = 1$, $C_{out} = 8$, kernel size = (1,2), stride = 1, dilation = 2], [$C_{in} = 8$, $C_{out} = 16$, kernel size = (1,3), stride = 1, dilation = 2] and [$C_{in} = 16$, $C_{out} = 32$, kernel size = (1,3), stride = 1, dilation = 2] respectively. Each batched output is normalized using the Batch Norm2d layer. The hidden dimension for the node feature project is set at 128. For optimization,

Table 5: Taxonomy of Graph Learning Methods for MTS Forecasting

Method	Learnable	Dynamic	Use Prior	Structure Form		Structure Availability	
				Binary	Continuous	Complete	Partial
STGCN	✗	-	✓	✗	✓	✓	✗
AGCRN	✓	✗	✗	-	-	-	-
MTGNN	✓	✗	✗	-	-	-	-
ESG	✓	✓	✗	-	-	-	-
StemGNN	✓	✓	✗	-	-	-	-
GTS	✓	✗	✓	✓	✗	✓	✗
SKI-CL	✓	✓	✓	✓	✓	✓	✓

we train SKI-CL with 100 epochs for every stage under Adam optimizer with a linear scheduler. For the learning rate schedule, we use a linear scheduler, which drops the linear rate from 0.0001 to the factor of 0.8 for every 20 epochs. The data split is 6/2/2 for training/validation/testing. We use a batch size of 32/64/64 for the Traffic-CL, Solar-CL and HAR-CL datasets and use a batch size of 8 for our synthetic dataset. Considering the sizes of datasets, the default memory for each regime is 1% for Traffic-CL and Synthetic-CL and 0.1% for Solar-CL and HAR-CL, respectively. We also weigh the examples in the current stage and in memory differently. We apply a weighted loss regarding the sizes of memory and training data, as stated in the manuscript, and we also construct a data loader that guarantees the balance between training data and memory data. For the setting of our distribution characterization scheme, the default values are $N = 10$ and $K = 7$. We implement our models in Pytorch. All experiments are run on one server with four NVIDIA RTX A6000 GPUs. **We will release our code upon paper acceptance.**

C.3 EVALUATIONS

Multi-horizon MTS Forecasting We use two common evaluation metrics for multi-horizon MTS forecasting, Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE), which are given as:

$$\text{MAE}(Y, \hat{Y}) = \frac{1}{\tau} \sum_{t=1}^{\tau} |y_t - \hat{y}_t| \quad (6)$$

$$\text{RMSE}(Y, \hat{Y}) = \sqrt{\frac{1}{\tau} \sum_{t=1}^{\tau} (y_t - \hat{y}_t)^2} \quad (7)$$

where τ is the number of time steps, \hat{y}_t is the forecasting results at t -th time step and y_t is the corresponding ground truth. Besides, \bar{y} and $\bar{\hat{y}}$ denote the mean values of ground truth and forecasting results, respectively.

Learning Faithful Dependency Structures For continuous edge variables, we still use MAE and RMSE to measure the similarity between the learned weighted graphs and continuous structural knowledge in an average sense, where the τ in Equations 6 and 7 denotes the entry of adjacency matrix. For binary edge variables, we use the average precision (Prec.) and recall (Rec.) to measure the similarity between the learned dependency structures over all testing MTS input windows and the binary structural prior at each regime, which are given as:

$$\text{Prec.} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (8)$$

$$\text{Rec.} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (9)$$

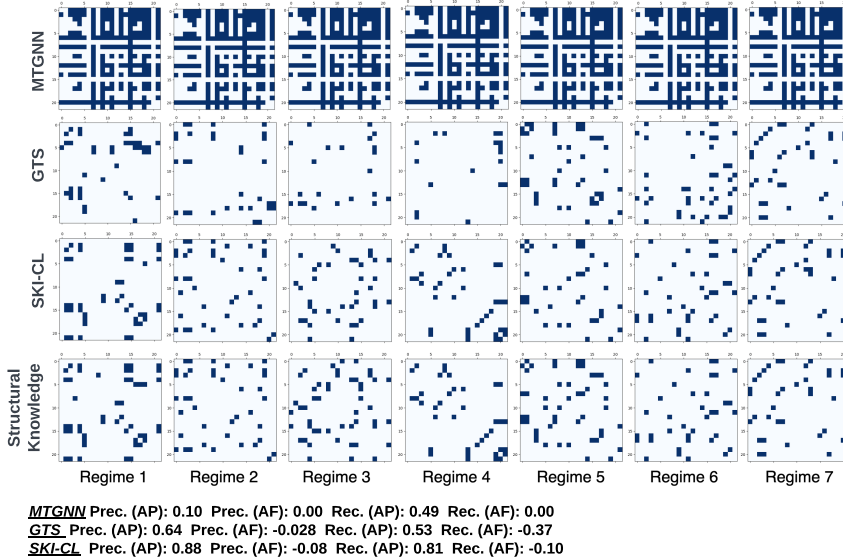


Figure 5: Visualization of Learned Structures on Traffic-CL Dataset.

where TP denotes the number of identified edges that exist in the structural prior, TN denotes the number of non-identified that do not exist in the structural prior, FP denotes the number of identified edges that do not exist in the structural prior, and FN denotes the number of non-identified edges that exist in the structural prior.

Continual MTS Forecasting and Dependency Structures Preserving We adopt two widely used metrics to evaluate the performance on continual MTS forecasting and dependency structures preserving, *i.e.*, the Average Performance (AP) and Average Forgetting (AF) Lopez-Paz & Ranzato (2017); Zhang et al. (2022a), where the AP and AF at i -th regime are defined as:

$$AP = \sum_{j=1}^i \frac{P_{i,j}}{i}, \forall i \geq 1 \tag{10}$$

$$AF = \sum_{j=1}^{i-1} \frac{(P_{i,j} - P_{j,j})}{i-1}, \forall i \geq 2 \tag{11}$$

where $P_{i,j}$ denotes the performance on regime j (including the forecasting performance and structure similarity) after the model has been sequentially trained from stage 1 to i .

Even if we provide both metrics for performance evaluation, we need to emphasize the superiority of average performance over average forgetting in continual learning. Average performance provides a direct measure of how well a learning system is performing on a task or set of tasks. It reflects the system’s ability to retain previously learned knowledge over past regimes while adapting to new information. While average forgetting is a relevant metric in assessing the memory capabilities of a learning system, it does not provide a complete picture of the learning system’s retention abilities. The average performance takes into account both the retention of old knowledge and the acquisition of new knowledge, providing a more comprehensive evaluation of the learning system’s performance. Therefore, we use average performance as the main evaluation metric and average forgetting as an auxiliary metric to measure knowledge retention and model adaptivity.

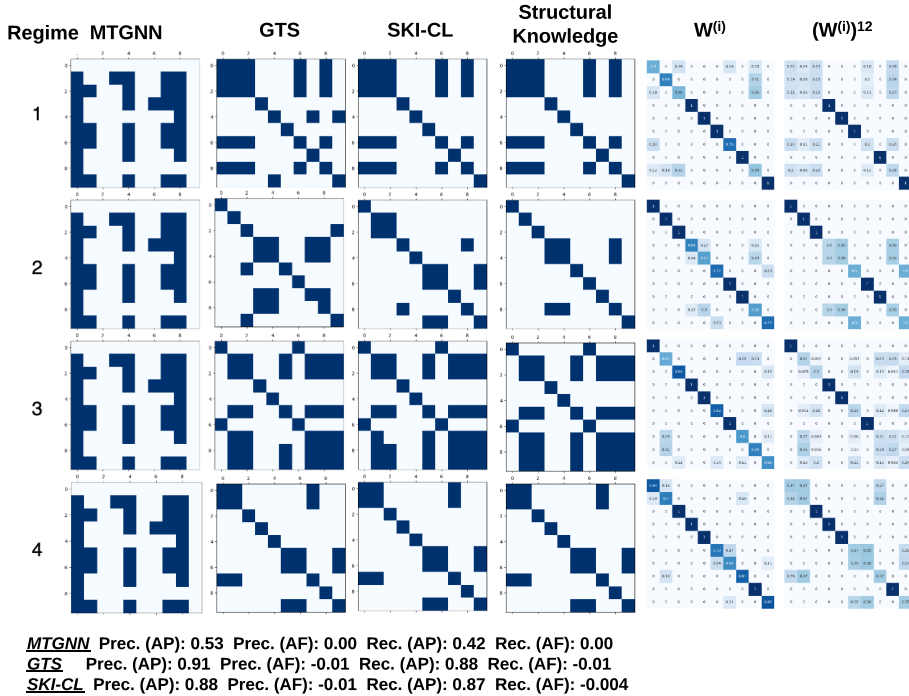


Figure 6: Visualization of Learned Structures on Synthetic-CL Dataset.

D VISUALIZATION OF LEARNED DEPENDENCY STRUCTURES

In this section, we provide case studies of the learned dependency structures for continual MTS forecasting on all datasets, as shown in Figure 5,6,7, and 8 . For binary edge scenarios, we compare our proposed SKI-CL with $MTGNN_{er}$ and GTS_{er} , which only support discrete edge formulation and yield better performance over other baselines. For continuous edge scenarios, we compare our SKI-CL with ESG_{er} that only support continuous edge formulation. For ESG_{er} , we visualize the generated graph at the last step representing the temporal dynamics of the whole sequence. All the results are plotted after the model is trained at the last regime, where the average structure similarities over all testing windows are also annotated under the case visualizations.

We first discuss the results on the Traffic-CL dataset as shown in Figure 5. As MTGNN directly learns a parameterized graph, it can only infer one fixed dependency structure reflecting the model at the latest regime. That being said, even if the model with experience replay is able to maintain the forecasting performance over the past regimes, the inferred graph fails to preserve the learned unique variable dependencies for each regime. Similarly, GTS infers a graph based on parameterized edge distributions, which also reveals the dependencies in a static sense. The incorporation of structural knowledge helps GTS characterize each regime in continual learning, which renders more relevant dependency structure for each regime. Nevertheless, learning a static edge distribution over regimes falls short of handling varying relational and temporal dynamics in complex environments. Another inherent drawback of GTS is that *GTS has to access the training data and memory buffer when inferring graphs at the testing stage, which is not realistic for practical model deployment in real-world applications*. Compared to these baselines, our model yields better forecasting performance (as shown in the manuscripts, as well as Table 8, and a faithful dependency structure that is more consistent with structure knowledge for each regime. These observations demonstrate the importance of the joint design of the dynamic graph inference module and the regularizer based on a structural prior.

For results on Synthetic-CL dataset as shown in Figure 6, the aforementioned observations still hold despite the comparative performance of GTS due to the simplicity of this dataset. Moreover, our model and GTS regularized with strong absolute correlations are able to render binary structures that

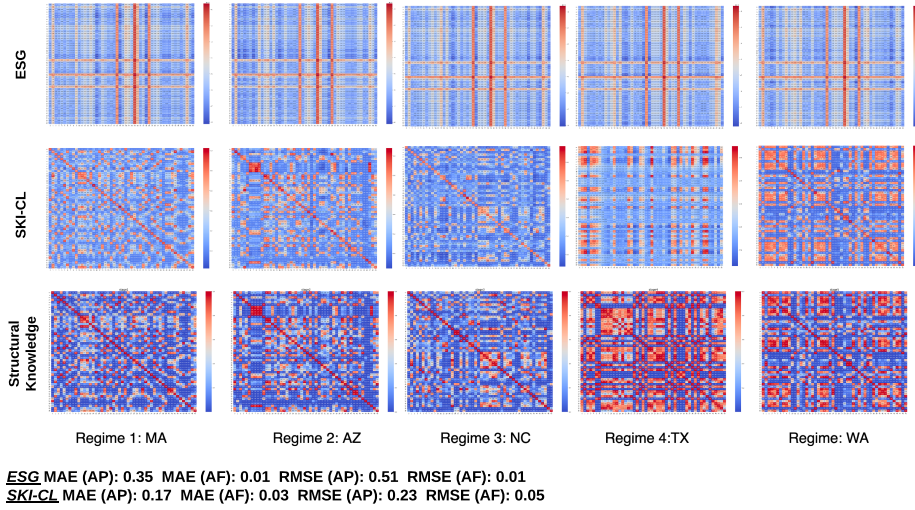


Figure 7: Visualization of Learned Structures on Solar-CL Dataset.

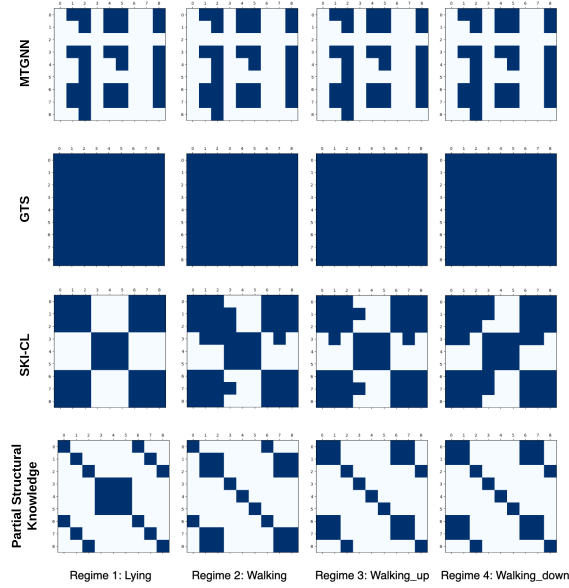


Figure 8: Visualization of Learned Structures on HAR-CL Dataset.

Table 6: Hyperparameter analysis for distribution characterization

N	K	Forecasting Performance($\times 10^{-2}$)				Structure Inference Performance			
		AP		AF		AP		AF	
		MAE	RMSE	MAE	RMSE	Prec.	Recall	Prec.	Recall
5	3	3.54	4.35	0.30	0.38	0.69	0.60	-0.11	-0.17
10	3	3.34	4.27	0.20	0.27	0.65	0.69	-0.04	-0.11
10	5	3.27	4.25	0.18	0.25	0.80	0.78	-0.01	-0.06
10	7	3.24	4.24	0.15	0.23	0.76	0.76	-0.06	-0.08
15	5	3.45	4.18	0.27	0.33	0.75	0.71	-0.06	-0.12
15	10	3.26	4.25	0.15	0.21	0.86	0.75	-0.06	-0.09

capture the single-step and multi-step variable interactions of the ‘ground truth’ dynamic adjacency matrix W^i and its 12th power (the input window size) for each regime.

We next discuss the results on the Solar-CL dataset as shown in Figure 7, where the edge is formulated as a continuous variable. It is clear that our proposed SKI-CL still gains advantages in terms of preserving a faithful continuous dependency structure for each regime. Instead, ESG that learns dynamic graphs still falls short of capturing a relevant structure due to the lack of regime characterizations.

Finally, we investigate our method on the HAR-CL dataset (as shown in Figure 8) when the structural knowledge is partially observed. Here, we do not evaluate the structure similarity due to the incompleteness of the prior as a referencing graph. Instead, we focus on how our model leverages the limited but confident knowledge in dependency structure learning. It can be seen that MTGNN fails to capture the important relationships that reveal in the structural knowledge. Besides, GTS consistently inferred a fully connected graph at the testing stages (even if we have tuned the temperature in the Gumbel-Softmax module), which renders a less meaningful dependency structure. Instead, the proposed SKI-CL exploits partial knowledge and renders faithful structures. For example, SKI-CL is able to capture the correlations of linear accelerations, angular velocities, and total accelerations within three axes, and the irrelevance between accelerations and angular velocities, which is reasonable in a binary sense for lying down behavior. Besides, even if the partial structural knowledge is the same for walking upstairs and walking downstairs, the SKI-CL is able to identify different dependencies patterns for different activities. It demonstrates the effectiveness of the SKI-CL on partially observed structural knowledge, suggesting a certain generalizability of our proposed framework in MTS modeling.

E HYPERPARAMETERS ANALYSIS

In this section, we supplement additional analysis of distribution characterization hyperparameters, namely the granularity N (10 by default) and the mode number K (7 by default), using Synthetic-CL dataset. As shown in Table 6, for a fixed N , a relatively large K facilitates inferred structure-preserving and mitigated the forgetting in time series forecasting. When N and K are close, average performance and average forgetting behavior are insensitive to the choice of these hyperparameters. However, a small N degrades the structure-preserving ability as the average forgetting on precision and recall increases.

F EXPERIMENTAL RESULTS FOR DIFFERENT HORIZON FORECASTING

Table 8 summarize the experiment results of baselines and our proposed SKI-CL method with variants for different horizon forecasting performance based on three rounds of experiments. We intentionally select 3-horizon prediction on Traffic-CL and Solar-CL datasets as the settings in Wu et al. (2020) and Cao et al. (2020). While for HAR-CL and synthetic-CL dataset, 6-horizon prediction performance is reported. Based on the results, the memory-replay-based methods generally alleviate the forgetting issues with better APs and smaller relative AFs. Under different horizon prediction settings, SKI-CL and its variants (SKI-CL_{er} and SKI-CL_{der++}) still consistently achieve the best or the second-best APs, demonstrating its advantages over other baseline methods.

Table 8: Experiment Results for Different Horizon Prediction (3-step Horizon Prediction for Traffic-CL and Solar-CL and 6-step Horizon Prediction for HAR-CL and Synthetic-CL) (Lower MAE and RMSE for AP mean better; When AP is comparable, lower MAE and RMSE for AF mean better.)

Model	Traffic-CL (3)				Solar-CL (3)				HAR-CL (6) ($\times 10^{-2}$)				Synthetic-CL (6) ($\times 10^{-2}$)			
	AP ↓		AF		AP ↓		AF		AP ↓		AF		AP ↓		AF	
	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
VAR _{seq}	73.36	104.97	70.26	97.19	59.30	196.75	72.88	242.41	18.04	27.37	1.01	1.28	18.61	27.33	7.70	11.40
ARIMA _{seq}	141.53	159.81	77.65	78.19	14.93	19.12	5.84	3.71	41.02	53.11	2.42	2.11	42.25	43.53	13.41	13.01
TCN _{seq}	13.89	22.12	2.57	4.77	0.93	2.71	0.12	0.13	12.71	20.46	2.98	4.35	3.81	3.37	0.26	0.33
TCN _{herd}	13.62	21.97	1.63	1.89	0.90	2.64	0.04	0.03	11.85	19.34	1.70	2.53	2.77	3.31	0.23	0.31
TCN _{er}	13.21	21.85	1.58	1.63	0.92	2.64	0.03	0.02	11.68	18.98	1.54	2.29	2.67	3.18	0.18	0.22
TCN _{der++}	13.50	21.90	1.48	1.57	0.87	2.63	0.03	0.05	11.63	18.97	1.46	2.11	2.73	3.29	0.15	0.17
Autoformer _{seq}	18.19	30.22	29.37	4.23	3.19	6.47	0.82	1.69	16.79	24.78	1.83	1.98	3.83	5.38	0.35	0.24
Autoformer _{herd}	17.72	2.67	1.89	3.02	2.89	5.34	0.67	1.23	16.01	24.12	1.58	1.78	3.71	5.26	0.31	0.20
Autoformer _{er}	16.07	26.82	1.36	2.54	2.54	5.02	0.13	0.23	15.66	23.88	1.57	1.58	3.65	5.26	0.24	0.17
Autoformer _{der++}	16.17	27.01	1.74	2.93	2.47	4.79	0.12	0.33	16.14	24.28	1.59	1.56	3.52	5.01	0.17	0.13
LSTNet _{seq}	24.89	39.49	10.50	15.63	2.45	4.77	1.09	1.45	14.13	21.90	4.05	5.49	25.31	33.71	7.44	10.61
LSTNet _{herd}	19.11	30.68	3.22	4.73	2.26	4.43	0.81	0.89	13.12	20.66	2.64	3.66	19.75	26.28	0.60	0.90
LSTNet _{er}	18.29	29.42	2.55	3.95	1.87	4.10	0.37	0.55	12.65	20.06	2.08	2.94	19.31	24.76	0.31	0.77
LSTNet _{der++}	17.94	28.96	2.03	3.18	1.78	3.98	0.18	0.28	12.55	19.94	1.85	2.86	18.92	24.26	0.77	0.57
STGCN _{seq}	26.80	42.22	11.90	14.48	3.14	6.30	1.86	3.16	16.52	25.04	5.21	6.52	8.21	12.90	2.27	4.63
STGCN _{herd}	26.53	43.56	9.36	14.79	2.57	5.99	1.15	2.65	15.41	23.97	4.12	5.69	5.53	8.65	1.42	3.37
STGCN _{er}	25.54	42.43	8.29	11.71	2.41	4.88	1.10	1.52	15.11	23.54	3.72	4.25	5.12	8.49	1.15	1.76
STGCN _{der++}	25.53	42.43	8.12	10.03	2.26	4.41	0.66	0.63	15.51	23.64	3.52	4.13	6.95	11.12	1.86	2.89
AGCRN _{seq}	16.13	26.21	3.65	5.29	1.44	3.21	0.31	0.60	13.82	21.88	2.97	4.06	13.00	21.16	0.99	0.47
AGCRN _{herd}	15.17	25.32	1.59	2.71	1.30	3.14	0.11	0.33	12.32	20.63	2.03	3.96	12.18	19.27	0.12	0.16
AGCRN _{er}	14.88	24.72	1.47	2.43	1.27	3.04	0.10	0.04	12.28	20.54	1.98	3.03	10.84	16.82	0.10	0.17
AGCRN _{der++}	15.02	25.15	1.41	2.66	1.17	3.01	0.03	0.03	12.01	20.18	1.76	2.95	12.40	19.56	0.25	0.23
StemGNN _{seq}	14.10	23.51	1.92	1.63	1.30	3.10	0.24	0.28	17.31	25.43	6.31	7.23	11.71	17.77	1.54	1.01
StemGNN _{herd}	13.78	23.01	0.98	1.03	1.17	2.92	0.22	0.25	16.51	25.03	5.21	4.23	9.17	13.97	1.03	0.65
StemGNN _{er}	13.49	22.59	0.36	0.42	1.00	2.73	0.04	0.05	16.43	24.83	5.19	4.13	8.38	12.95	0.96	0.97
StemGNN _{der++}	13.37	22.42	0.31	0.37	1.03	2.87	0.04	0.07	16.25	24.51	5.10	4.01	9.02	13.65	0.89	0.85
ESG _{seq}	14.62	23.83	3.72	5.39	1.75	2.93	0.92	0.83	16.57	25.07	5.34	6.21	8.68	13.95	1.06	1.57
ESG _{herd}	14.53	23.79	3.26	4.71	1.01	2.78	0.14	0.27	14.31	23.17	4.48	5.39	7.61	12.25	0.71	1.15
ESG _{er}	13.07	22.11	1.62	2.93	1.15	2.71	0.32	0.25	14.21	23.01	4.37	5.21	7.21	11.65	0.61	0.89
ESG _{der++}	13.92	23.17	2.52	3.98	0.98	2.67	0.16	0.21	14.10	22.87	4.13	4.83	8.08	12.95	0.35	0.52
GTS _{seq}	15.20	27.09	3.23	7.32	1.20	3.19	0.19	0.48	14.44	23.41	3.51	4.63	4.90	6.71	0.56	0.61
GTS _{herd}	14.33	23.55	2.82	3.78	1.03	2.69	0.05	0.10	13.78	22.89	3.08	4.05	4.71	6.66	0.13	0.20
GTS _{er}	14.54	24.65	1.96	3.58	0.94	2.67	0.08	0.05	13.38	22.45	2.95	3.97	4.41	6.40	0.13	0.19
GTS _{der++}	14.07	24.12	1.25	3.39	0.96	2.67	0.05	0.15	13.34	22.41	2.93	3.89	4.39	6.39	0.12	0.18
MTGNN _{seq}	14.42	23.94	3.34	5.13	1.14	2.76	0.32	0.30	12.87	20.66	3.51	4.80	7.63	12.44	0.59	0.67
MTGNN _{herd}	13.55	22.55	2.13	3.35	0.96	2.69	0.13	0.23	11.98	19.78	2.42	3.61	7.43	11.99	0.45	0.61
MTGNN _{er}	13.95	22.29	2.03	3.02	0.95	2.74	0.05	0.16	11.30	18.67	1.51	2.15	6.38	10.13	0.35	0.53
MTGNN _{der++}	12.99	21.86	1.24	2.24	0.91	2.72	0.03	0.10	11.27	18.64	1.22	1.81	6.31	10.10	0.18	0.41
SKI-CL _{seq}	13.81	23.16	2.52	4.25	0.90	2.63	0.09	0.14	12.41	19.37	2.54	3.78	3.72	5.26	2.02	3.10
SKI-CL _{herd}	12.91	21.80	1.43	2.55	0.87	2.65	0.07	0.12	11.78	19.04	1.89	2.95	3.11	4.18	0.23	0.45
SKI-CL _{er}	12.66	21.21	1.26	2.27	0.88	2.64	0.05	0.06	11.67	18.95	1.54	2.35	2.85	3.60	0.17	0.16
SKI-CL _{der++}	12.61	21.20	1.10	1.72	0.85	2.56	0.04	0.05	11.55	18.84	1.40	2.06	2.33	3.05	0.10	0.13
SKI-CL	12.49	21.01	0.93	1.60	0.82	2.54	0.02	0.02	11.01	18.25	1.21	1.78	2.01	2.85	0.05	0.03

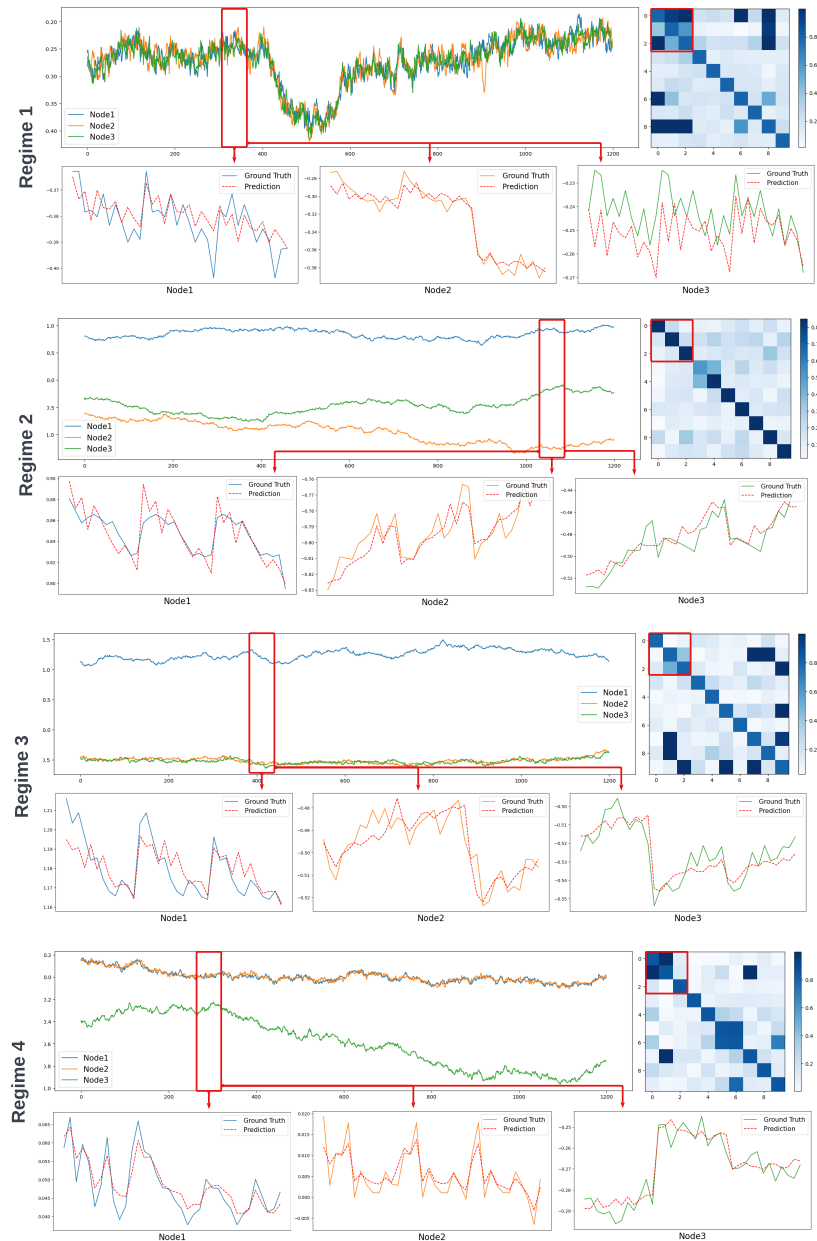


Figure 9: A case study of SKI-CL on Synthetic-CL dataset. In each regime, red rectangles indicate the correspondence between the ground truth time series (top-left) and the inferred variable dependencies (top-right), red arrows indicate the comparisons between these ground truth values and corresponding predictions (bottoms).

G A CASE STUDY WITH INFERRED DEPENDENCY STRUCTURES AND PREDICTION VISUALIZATIONS

We provide a case study on the Synthetic-CL dataset to illustrate the efficacy of SKI-CL, as shown in Figure 9. Our analysis is based on the final SKI-CL model that has been sequentially trained over all regimes. We select three variables (nodes) and visualize the testing data, where the temporal dynamics obviously differ across four regimes. It is clear that SKI-CL can render a faithful dependency structure that well aligns the similarity of variables in each regime (*e.g.*, three nodes are highly similar to each other in regime 1; only node 1 and node 2 are similar in regime 4). Moreover, SKI-CL gives relatively accurate forecasts that capture each variable’s temporal dynamics of ground truths.