

CAST-GNN: Continual Adaptive Learning for Custom Spatio-Temporal Knowledge Graphs via Graph Neural Networks

Gözde Ayşe Tataroğlu Özbek*[†], Yash Raj Shrestha*, Jean-Paul Calbimonte^{†‡}

* University of Lausanne, Lausanne, Switzerland

[†] University of Applied Sciences and Arts Western Switzerland HES-SO, Sierre, Switzerland

[‡] The Sense Innovation and Research Center, Lausanne, Switzerland

Emails: gozdeayse.tatarogluzbulak@unil.ch; yashraj.shrestha@unil.ch; jean-paul.calbimonte@hevs.ch

Abstract—Real-time video streams present unique challenges for continual learning systems, demanding models that can incrementally update representations, preserve past knowledge, and reason over complex semantic relationships without sacrificing efficiency. In this paper, we introduce CAST-GNN, the first unified Graph Neural Network architecture expressly designed for continual adaptation on streaming Spatio-Temporal Knowledge Graphs (STKGs) derived from open-source video benchmarks. CAST-GNN integrates dynamic temporal embedding layers, adaptive self-attention, episodic graph pattern memory, and a novel hybrid selective replay buffer with Fisher-based regularization and knowledge distillation to mitigate catastrophic forgetting. Through comprehensive experiments on four diverse STKG benchmarks (UCF-101, HMDB-51, Kinetics-400, and Something-Something), our model achieves 96–97% accuracy, between 0.13–0.31% forgetting, by consistently outperforming re-implemented continual-learning baselines under identical conditions. Ablation studies confirm the critical synergy between temporal embeddings and adaptive attention. We further demonstrate XAI-driven interpretability by aligning global distributional shifts with local node-level attributions. CAST-GNN not only advances robust semantic reasoning and knowledge retention but also provides a scalable, explainable framework applicable to a wide array of real-world streaming scenarios.

Index Terms—Continual Learning, Spatio-Temporal Knowledge Graphs, Graph Neural Networks, Catastrophic Forgetting, Stream Reasoning

I. INTRODUCTION

Recent advancements in real-time streaming data analysis have emphasized the necessity for dynamic, continually adaptive machine learning models, particularly in applications involving video understanding, autonomous driving, or sensor network management. With the exponential growth of streaming data sources, effectively modeling temporal and spatial dynamics within continuous learning frameworks becomes increasingly essential yet highly challenging. Traditional Graph Neural Networks (GNNs) have demonstrated remarkable capabilities in modeling relational structures within static datasets. In this context, Spatio-Temporal Knowledge Graphs (STKGs), which structure visual and temporal data into relational representations, have emerged as a promising foundation for capturing and reasoning over complex semantic relationships among objects and actions. Despite significant progress, exist-

ing methods exhibit critical limitations such as integrating real-time adaptability with robust semantic reasoning. Also, most of methods are struggling to mitigate catastrophic forgetting effectively in dynamic and streaming environments.

Indeed, traditional continual learning approaches tend to isolate structural adaptations from semantic-level interpretation, leading to fragmented solutions that either excel in incremental embedding updates or in preserving long-term knowledge, but rarely both. Furthermore, current frameworks typically lack mechanisms for interpreting model decisions, limiting their applicability in real-world scenarios where explainability is crucial. There remains a significant gap in developing unified, scalable solutions capable of continuously adapting to new information while preserving semantic coherence and preventing catastrophic forgetting.

Motivated by these limitations, this paper introduces CAST-GNN model, a novel Spatio-Temporal Graph Neural Network specifically designed for continual learning scenarios on streaming STKGs derived from video data. Our model uniquely integrates temporal embedding techniques, adaptive self-attention, episodic memory structures, selective replay buffers, and advanced regularization strategies within a unified, coherent architecture. To the best of our knowledge, CAST-GNN represents the first holistic approach that simultaneously addresses dynamic embedding adaptability, catastrophic forgetting mitigation, semantic reasoning, and interpretability for real-time STKG-based continual learning.

The primary contributions of our work are the following:

- *Unified Continual Learning Framework for Streaming STKGs:* We introduce CAST-GNN, the first end-to-end continual learning architecture explicitly designed for spatio-temporal knowledge graphs derived from real-time video streams by unifying structural adaptability with rich semantic interpretability.
- *Robust Semantic Reasoning with Contextual Coherence:* Our method integrates adaptive self-attention and temporal embedding layers to maintain high-level semantic consistency even under unexpected shifts owing to semantic-level reasoning capabilities in streaming environments.

- *Advanced Catastrophic Forgetting Mitigation:* Our model structure combines Fisher information-based regularization and knowledge distillation with a novel hybrid selective replay buffer. It achieves optimal balance between rapid adaptation and long-term knowledge retention.
- *Adaptability and Explainability for Trustworthy Deployment:* Our model not only adapts swiftly to diverse, dynamically changing data distributions, but also provides transparent, node level explanations by enhancing user trust and facilitating deployment across domains.
- *Generalizability and Practical Scalability:* Extensive experiments on multiple benchmarks (from activity recognition in video analytics to urban mobility forecasting and sensor network analysis) demonstrate our model’s robust performance and scalability by underscoring its potential for broad real-world application.

In summary, our proposed model, CAST-GNN, represents a significant leap forward in continual graph learning methodologies, effectively integrating real-time adaptability, semantic-level reasoning, and robust forgetting mitigation into a unified, scalable solution. Through comprehensive experiments and evaluations, we validate its efficacy and versatility by positioning CAST-GNN as a foundational framework for future research and practical deployments in dynamic, real-world applications.

II. RELATED WORK

Dynamic Graph Representation Learning: Early efforts in dynamic graph learning focused on efficiently updating node embeddings as new edges or events arrived. Approaches such as Streaming Graph Neural Networks (DyGNN) [1], Temporal Graph Networks (TGN) [2], DySAT [3], and CTDNE [4] have introduced a range of mechanisms from recurrent parameter evolution to attention-based temporal encoding for capturing evolving graph structures. While these models enable online adaptation to changing topologies, they generally do not address the compounded challenges of catastrophic forgetting in long-running streams. In addition, they do not explicitly support high-level semantic reasoning over STKG data as in our case.

Mitigating Catastrophic Forgetting: In parallel, continual learning research has developed replay and regularization strategies to preserve previously acquired knowledge. For instance, the TIE framework applies experience replay and temporal regularization for incremental knowledge graph completion [5], and History Repeats extends this strategy with clustered replay for event-centric temporal KGs [6]. Regularization-based methods such as Elastic Weight Consolidation (EWC) [7] leverage Fisher information to penalize significant updates to crucial parameters. Recent surveys in continual graph learning [8], [9] confirm that while these approaches alleviate forgetting, they often remain disconnected from dynamic embedding updates and are rarely evaluated in real-time, video-derived STKG scenarios.

Adaptive Meta-Learning and Semantic Reasoning: Meta-learning has recently been applied to improve adaptation under distributional shifts. TGOonline employs online meta-learning for temporal graphs [10], and DOST proposes distribution-aware continual learning for urban spatio-temporal forecasting [11]. Despite these advances, these methods typically focus on prediction accuracy in non-stationary environments, rather than on semantic activity recognition or explicit reasoning over knowledge graphs as in our case. Until now, no prior work provides a unified GNN framework that simultaneously incorporates incremental embedding, advanced forgetting mitigation like replay and Fisher regularization, and semantic-level reasoning for real-time streaming in video-derived STKGs.

While the previously mentioned methods have advanced the state of dynamic graph learning and continual adaptation, none effectively tackle the simultaneous need for online embedding updates, semantic-level reasoning, and robustness to catastrophic forgetting in streaming video-derived STKGs. In contrast, our proposed CAST-GNN framework integrates state-of-art modules in a unified architecture. This unique approach addresses the fragmentation present in the current literature (advancing beyond isolated online updates, distinct forgetting mitigation strategies, or limited semantic reasoning) to provide a holistic solution that is designed for real-time continual learning and semantic reasoning on STKGs.

III. DATASET

In this work, we address the longstanding gap in modeling the rich semantic relationships inherent in streaming visual data by constructing Spatio-Temporal Knowledge Graphs (STKGs) from publicly available video benchmarks. We first pre-process open-source datasets—each of which offers diverse scenes, object annotations, and standardized evaluation protocols—to convert raw video frames into graphs whose nodes represent detected objects and whose edges capture their spatio-temporal interactions. This transformation bridges the divide between unstructured visual streams and structured semantic representations, enabling continual learning methods to operate directly on the evolving relational patterns of real events. By using openly shared video datasets, our approach guarantees reproducibility and facilitates careful comparison with prior work in video understanding and dynamic graph learning. Through STKGs, we preserve both spatial and temporal contextual links, thereby overcoming the information loss and interpretability challenges that arise when semantics are compressed or ignored, and delivering a unified framework for real-time continual learning and semantic inference.

We evaluate CAST-GNN on four publicly available video benchmarks: HMDB-51 [12], UCF-101 [13], Kinetics-400 [14], and Something-Something V2 [15]. We uniformly preprocess frames and release the custom STKG datasets together with the source code of our proposed model, all publicly available at our Git repository¹. It should also be

¹<https://github.com/aislab-hevs/CAST-GNN>

noted that this repository also contains our main model’s full code and the datasets. It also provides documentation explaining how we created these custom STKGs from open source video-based datasets, for reproducibility and adaptability to different domains. We processed 120 videos from each dataset and created 15-class STKGs. Details of this custom STKGs construction process are provided in [16]. The number of classes was not kept too high so that the performance of the model could be observed more easily. Also, we present the main characteristics and differences of these STKGs in Table I, using the same names of the open-source dataset from which they are derived. These datasets show different distributions and diversity, as described in Table I, and for which our framework’s robustness is intended to be proven.

IV. PROPOSED MODEL

We propose CAST-GNN, an innovative Spatio-Temporal Graph Neural Network framework specialized specifically for incremental and continual learning on STKGs derived from real-time streaming video data. Our model addresses four critical challenges inherent in STKG modeling: incremental embedding updates, catastrophic forgetting mitigation, rapid adaptation to distributional changes, and efficient semantic-level reasoning.

A. Dynamic Spatio-Temporal Embedding Module

To effectively model the temporal dynamics and spatial relationships among nodes, our proposed model, CAST-GNN, introduces a temporal embedding strategy that jointly incorporates structural and temporal information. Specifically, the input node features \mathbf{x} and normalized timestamps \mathbf{t} are transformed as follows: structural representations are obtained via a Graph Convolutional Network (GCN), and temporal embeddings are learned through a linear transformation followed by a ReLU activation:

$$\mathbf{h}^{(1)} = \text{ReLU}(\text{GCN}(\mathbf{x}, \mathbf{E})), \quad \mathbf{e}_t = \text{ReLU}(\mathbf{W}_t \mathbf{t} + \mathbf{b}_t) \quad (1)$$

Here, $\mathbf{h}^{(1)}$ denotes the node embeddings obtained after the first GCN layer, while $\mathbf{h}^{(0)}$ corresponds to the raw input features. The structural embeddings $\mathbf{h}^{(1)}$ and temporal embeddings \mathbf{e}_t are concatenated and passed through a learnable transformation layer with dropout to enhance temporal awareness and improve model generalization:

$$\mathbf{h}_{\text{temp}} = \text{Dropout}(\mathbf{W}_f [\mathbf{h}^{(1)}; \mathbf{e}_t] + \mathbf{b}_f) \quad (2)$$

Here, \mathbf{W}_t , \mathbf{b}_t , \mathbf{W}_f , and \mathbf{b}_f are learnable parameters. ReLU is used as the non-linear activation function for both structural and temporal embeddings.

To further capture nuanced temporal interactions between connected nodes, we integrate a temporal self-attention mechanism inspired by DySAT [3]. This mechanism explicitly incorporates temporal differences Δt between nodes to allow adaptive modeling across varying temporal scales. In the equation below, \mathbf{E}_{attr} represents edge-level attributes and

$\mathbf{e}_{\Delta t}$ encodes temporal differences, enabling CAST-GNN to selectively emphasize temporally meaningful connections:

$$\mathbf{A}_{\text{temp}}(\mathbf{h}) = \text{TransformerConv}(\mathbf{h}, \mathbf{E}, [\mathbf{E}_{\text{attr}}; \mathbf{e}_{\Delta t}]) \quad (3)$$

B. Continual Learning and Adaptation Module

A major challenge in continual learning is catastrophic forgetting, where new information disrupts previously acquired knowledge. CAST-GNN addresses this issue by incorporating a set of complementary mechanisms aimed at preserving long-term memory while enabling rapid adaptation to new data.

To retain historical context, we employ an Episodic Graph Pattern Memory (EGPM), inspired by the temporal modeling approach in JODIE [17]. We also adopt the principles of Gradient Episodic Memory (GEM) [18], which emphasizes the importance of rehearsal-based mechanisms. In our case, a Gated Recurrent Unit (GRU) sequentially encodes temporal embeddings, enabling the model to update its internal memory states over time:

$$\mathbf{h}^{(2)}, \mathbf{h}_{\text{state}} = \text{GRU}(\mathbf{h}_{\text{temp}}, \mathbf{h}_{\text{state}}^{\text{prev}}) \quad (4)$$

In this formula, $\mathbf{h}^{(2)}$ denotes the updated node embeddings after temporal encoding, while $\mathbf{h}_{\text{state}}$ represents the recurrent hidden state that carries information across time steps. To further enhance temporal awareness, we incorporate Temporal Graph Network Memory (TGNMemory) [2], which continuously maintains node-specific memory vectors. At each timestamp t , the memory of node i is updated as:

$$\mathbf{m}_i^{(t)} \leftarrow \text{GRU}(\mathbf{m}_i^{(t-1)}, \mathbf{e}_i^{(t)}) \quad (5)$$

Here, $\mathbf{e}_i^{(t)}$ is the embedding of node i at time t , and $\mathbf{m}_i^{(t-1)}$ is its memory from the previous step. This recurrent update allows the model to capture temporal continuity in evolving graph streams.

For improving the adaptation under non-stationary distributions and concept drifts, we introduce a meta-learning module inspired by TGOOnline [10]. While lightweight in formulation, this block functions as an adaptive transformation layer; for consistency, we also refer to it as the *meta-learning module* in our ablation study:

$$\begin{aligned} \mathbf{u}_{\text{meta}} &= \text{LayerNorm}(\mathbf{W}_{\text{meta}} \mathbf{h}^{(2)} + \mathbf{b}_{\text{meta}}) \\ \mathbf{h}_{\text{meta}} &= \text{Dropout}(\text{ReLU}(\mathbf{u}_{\text{meta}})) \end{aligned} \quad (6)$$

In addition to memory-based mechanisms, we employ two regularization strategies to mitigate forgetting: Elastic Weight Consolidation (EWC) [7] and Knowledge Distillation (KD) [19].

EWC penalizes large changes to parameters that were important for previous tasks by weighting their deviation according to the Fisher information:

$$\mathcal{L}_{\text{EWC}} = \lambda \sum_i F_i \cdot (\theta_i - \theta_i^*)^2 \quad (7)$$

TABLE I: Characteristics of video datasets and their corresponding derived STKGs

Feature	UCF101	HMDB51	Kinetics400	Something-Something
Video Length	Short (5–10 s)	Very short (3–6 s)	Short–medium (~10 s)	Very short (2–6 s)
Total # of Classes	101	51	400	174
Activity Complexity	Medium	Low–medium	Medium–high	Medium–high (fine-grained)
Temporal Relationships	Medium	Low–medium	Medium–high	Very high
Content Type	Sports & daily life	Simple human motions	Mixed activities	Fine-grained object–action
Dataset Size (#video)	~13 000	~7 000	>300 000	>100 000
Example Activities	Bowling, Diving	Eating, Laughing	Playing Guitar, Skiing	Moving Something, Putting On
Node Count of Derived STKG	40301	33873	34241	7293
Edge Count of Derived STKG	366658	515396	222513	28598

where θ_i are current parameters, θ_i^* are previously learned parameters, and F_i represents the estimated importance of each parameter.

KD complements this by enforcing consistency between the current model (student) and its previous state (teacher). Instead of relying solely on ground-truth labels, the student also learns from the teacher’s probability distributions over all classes (i.e., softmax outputs), which carry valuable information about class relationships and previously learned behavior. The combined regularization loss is:

$$\mathcal{L}_{\text{reg}} = \alpha \cdot \text{KL}(p_{\text{teacher}} \parallel p_{\text{student}}) + \mathcal{L}_{\text{EWC}} \quad (8)$$

Here, p_{teacher} and p_{student} denote the probability distributions of teacher and student models, and α , λ are hyperparameters controlling the strength of the KD and EWC terms.

To support long-term knowledge retention, we integrate a Hybrid Selective Replay Buffer inspired by [20]. This buffer stores a curated set of representative past samples using a novel *uncertainty-class* sampling method, which combines uncertainty-based selection with class balancing. To maintain efficiency and diversity, the buffer is periodically refreshed via K-Medoids clustering [21]. The effectiveness of this strategy is demonstrated in our ablation study (Table III). During training, $\mathcal{L}_{\text{replay}}$ corresponds to the standard cross-entropy loss computed on samples drawn from the replay buffer \mathcal{B} .

The total training objective integrates the main classification loss with the proposed regularization and memory replay terms:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{reg}} + \mathcal{L}_{\text{replay}} \quad (9)$$

Here, \mathcal{L}_{cls} denotes the class-wise focal loss [22], which helps to address class imbalance issues commonly present in our real-time STKG data. $\mathcal{L}_{\text{replay}}$ represents the selective memory replay loss, which reinforces learning from representative samples of earlier data to further reduce forgetting. By jointly optimizing these loss components, our approach effectively preserves performance based on prior knowledge while adapting to new data in a continuous learning environment.

Through the careful integration and coordination of these specialized modules, our proposed CAST-GNN model addresses several critical gaps identified in prior literature [2], [6], [23]. These modules include dynamic temporal embedding layers, adaptive self-attention mechanisms, episodic graph pattern memory, temporal graph network memory, meta-learning

strategies, selective replay buffers, and sophisticated regularization methods. The coherent integration of these modules significantly enhances the semantic level reasoning capabilities while robustly reducing catastrophic forgetting, thus enabling efficient incremental learning. In conclusion, our model CAST-GNN effectively supports real-time continuous learning and successfully performs activity recognition on our custom STKG streaming data, while showing versatile applicability in various fields.

V. EXPERIMENTAL ANALYSIS

A. Evaluation Metrics

To comprehensively evaluate CAST-GNN’s performance, we utilize a combination of standard and continual learning-specific metrics to align with our research objectives. First, accuracy and F1-score are used to assess the overall classification capability by reflecting the balance between precision and recall. For deeper performance analysis, especially in ranking correct predictions among alternatives, we adopt Recall@5 and Mean Reciprocal Rank (MRR). These metrics highlight the model’s effectiveness in prioritizing relevant predictions, essential for real-time decision-making scenarios as in our case. In the context of continual learning, stability and robustness over time become critical. Therefore, we employ Temporal Coherence to explicitly measure prediction consistency across sequential inputs. To quantify the corruption of knowledge maintenance, we utilize the Windowed Forgetting Rate which provides a localized view of performance decline due to catastrophic forgetting within defined intervals. Furthermore, we report the Average Forgetting Rate to reflect the model’s long-term keeping of previously acquired knowledge [24]. Moreover, Adaptation Speed is used to measure how quickly our model adjusts to dynamic environments. While its value reduces, it validates the model’s quick adaptation capability. All formulas of these continual learning-based performance metrics are presented below:

$$\text{TC} = 1 - \frac{1}{N-1} \sum_{i=1}^{N-1} \mathbb{I}(\hat{y}_i \neq \hat{y}_{i+1}) \quad (10)$$

$$\text{WFR}_w(t) = \frac{1}{w} \sum_{j=t-w}^{t-1} \max(0, a_j - a_t) \quad (11)$$

$$\text{AFR} = \frac{1}{C} \sum_{c=1}^C \left(\max_e a_{c,e} - a_{c,E} \right) \quad (12)$$

TABLE II: Performance comparison of our method with similar approaches in the literature, by using our STKGs derived from four Open-Source Datasets

Methods	(%) Accuracy	(%) Forgetting Rate	Adaptation Speed (epochs)	Latency (ms)
UCF				
Our Method (Cast-GNN)	97.82	0.24	7.57	17.55
EvolveGCNN [23]	76	0.82	41.82	91.24
TGOnline [10]	87.87	7.15	9.00	3.13
DOST [11]	83	4.00	35.00	1616
History Repeats [6]	77	0.57	6.93	1.32
HMDB				
Our Method (Cast-GNN)	96.95	0.13	8.14	10.92
EvolveGCN [23]	74	0.79	43.06	23.25
TGOnline [10]	81	12.28	4.00	3.95
DOST [11]	89	1.17	22.00	1439
History Repeats [6]	75	0.55	7.06	1.02
Kinetics				
Our Method (Cast-GNN)	97	0.17	7.97	10.06
EvolveGCNN [23]	68	0.86	42.83	19.35
TGOnline [10]	82	7.91	11.00	9.69
DOST [11]	80.5	3.64	52.00	13.69
History Repeats [6]	65	0.47	7.56	1.01
Something-Something				
Our Method (Cast-GNN)	80	0.31	8.70	5.73
EvolveGCNN [23]	46	1.93	39.41	9.8
TGOnline [10]	70	12.12	8.00	1.98
DOST [11]	77	1.75	6.00	185
History Repeats [6]	58	0.51	8.98	0.81

$$\text{AdaptationSpeed} = \frac{1}{|\mathcal{C}'|} \sum_{c \in \mathcal{C}'} e_c \quad (13)$$

where:

- \hat{y}_i is the predicted label at time-step i , N is the total number of steps.
- a_t is the accuracy at time t , and w is the window size.
- For each chunk $c \in \{1, \dots, C\}$, $a_{c,e}$ is its accuracy at epoch e , and E is the final epoch.
- $\mathcal{C}' = \{c : e_c \text{ is defined}\}$ and e_c is the epoch when chunk c first reaches its maximum accuracy.

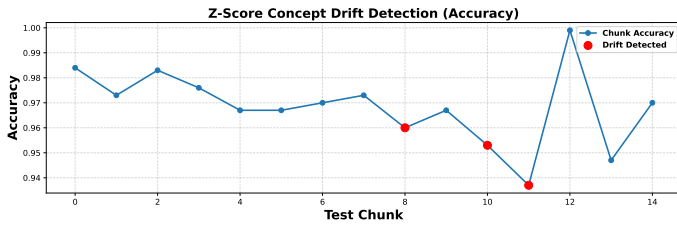
For clarity of notation, we abbreviate the key continual-learning metrics as follows: TC (Temporal Coherence), WFR (Windowed Forgetting Rate), AFR (Average Forgetting Rate), and AdaptationSpeed.

In addition to these metrics, we incorporate Z-score and ADWIN methods for concept drift detection to proactively identify and respond to changes in data distributions [25]. Thus we can observe our model coherence and timely adaptability in dynamically evolving streaming scenarios. Finally, we include Inference Latency to ensure the model’s suitability for real-time applications, since minimal latency is crucial for timely responses in streaming scenarios. In brief, these metrics were strategically chosen to robustly validate our model’s effectiveness in addressing continual learning challenges while providing applicability to real-time semantic reasoning tasks.

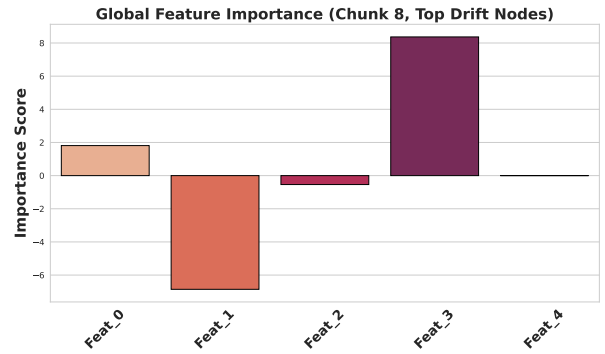
B. Comparative Analysis

We performed extensive experiments to evaluate the performance of our CAST-GNN method against closely related continual learning-based methods in dynamic and online scenarios. Although the domains of the compared approaches slightly differ from ours, they share common objectives, specifically continual adaptation and effective learning in dynamic settings. Importantly, since no existing studies in the literature explicitly address continual semantic reasoning tasks on spatio-temporal knowledge graphs (STKG) derived from open-source video datasets, we adapted related approaches. Specifically, EvolveGCN [23], TGOnline [10], DOST [11], and History Repeats [6] were adapted to our context and evaluated on STKGs derived from the UCF, HMDB, Kinetics, and Something-Something datasets. It should be emphasized that the datasets themselves were not directly used, but rather transformed into STKGs encapsulating semantic and spatio-temporal relationships suitable for reasoning tasks.

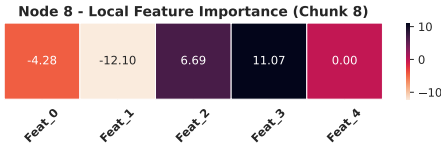
According to Table II, our CAST-GNN consistently outperforms all comparative methods across all datasets in terms of accuracy, average forgetting rate, adaptation speed, and latency. Particularly notable is our low forgetting rate and rapid adaptation speed, demonstrating robustness in incremental learning tasks and quick responsiveness to distributional shifts. For instance, CAST-GNN achieves exceptionally low forgetting rates (0.13%-0.31%) compared to other methods such as TGOnline (7.15%–12.28%) and EvolveGCN (0.79%–1.93%), while DOST stays around 1.17%–4.00%. These ranges highlight our model’s superior stability in retain-



(a) Concept drift detection over chunks using Z-Score.



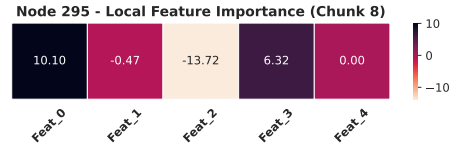
(b) Global feature importance (Chunk 8, Top Drift Nodes).



(c) Node 8 - Local importance.



(d) Node 23 - Local importance.



(e) Node 295 - Local importance.

Fig. 1: XAI visualization for the chunk where the first drift is detected. (a) Drift detection via accuracy drop. (b) Global feature importance in the drifted chunk. (c-e) Local feature importance for top-3 drifted nodes. This figure shows how CAST-GNN detects concept drift and explains it by linking global feature distributions with local node-level attributions. The aim is to uncover which features drive these shifts and how they are reflected at both global and local levels.

ing knowledge over time. In terms of specific datasets, our approach achieves outstanding performance on UCF and Kinetics datasets, characterized by clearly defined action semantics and structured temporal patterns, attaining highest accuracy (97.82% and 97% respectively) and minimal forgetting (0.24% and 0.17% respectively). Interestingly, all methods, including ours, encounter lower accuracy and higher forgetting on the Something-Something dataset. This observation underscores the dataset’s inherent complexity, subtle temporal dependencies, and different semantic differences, which pose substantial challenges for continual learning. Nevertheless, CAST-GNN maintains a relatively superior performance even on this challenging dataset, affirming its effectiveness in handling complex semantic dynamics. Comparatively, EvolveGCN demonstrates higher latency and slower adaptation speeds across all datasets, attributed to its confidence on evolving graph structures and recurrent-based architectures, while TGOOnline, benefiting from adaptive meta-learning, shows better adaptation speeds but suffers from higher forgetting rates, especially visible on HMDB and Something-Something datasets. History Repeats provides minimal inference latency due to its event-centric knowledge graph completion approach but does not match CAST-GNN’s accuracy and robustness, reflecting its suitability mainly for event-centric rather than activity recognition tasks. Finally, DOST, designed for urban spatio-temporal forecasting, unsurprisingly shows high latency and adaptation speeds unsuited for real-time video-based scenarios. In summary, our CAST-GNN demonstrates superior and balanced performance across multiple metrics and diverse datasets, significantly advancing the state-of-the-art for continual semantic reasoning in STKG-

based real-time activity recognition scenarios.

Figure 1 illustrates our explainability pipeline in continual GNN learning through a spatio-temporal graph setting, structured from top to bottom in three complementary stages. At the top (a), the Z-Score concept drift detector highlights the specific chunks (red markers) in the test accuracy curve where significant distributional changes occur. The very first flagged chunk represents the point at which the model’s learned representation begins to misalign with the incoming data distribution. Thereby it signals the onset of drift and motivating further inspection. In the middle (b), a global feature importance bar chart is computed via Integrated Gradients (GNNShap) [26]. This captures the aggregated contribution of features across the top three nodes which exhibit the highest drift within that critical chunk. This global view provides a bridge between the detected distributional shift and the features most responsible for it. Therefore a clear cause-and-effect link between overall accuracy drops and node-level dynamics are established. At the bottom (c–e), three heatmaps provide local (node-specific) feature importance scores for the very same high-drift nodes (Node 8, Node 23, and Node 295). Notably, the fine-grained patterns in each heatmap closely mirror the peaks and distributions observed in the global feature profile. This pattern confirms that the global drift is underpinned by node-level feature behavior. Also, this multi-level explainability pipeline is novel in the context of continual learning on spatio-temporal knowledge graphs. The pipeline detects concept drift, then aligns it with both global feature distributions and local node-level attributions. Rather than treating nodes in isolation, CAST-GNN’s predictions emerge

from inter-node interactions and temporal context encoded in our STKG. The global analysis captures collective feature shifts across the graph, while the local attributions reveal how individual nodes recalibrate their feature importance when confronted with distributional change. By aligning these complementary perspectives, we demonstrate that CAST-GNN meaningfully leverages semantic relationships in streaming video-derived graphs and enables each node to adjust its feature confidence in real time as new data arrives.

Figure 2 provides a comprehensive analysis of the temporal coherence and windowed forgetting trends of our CAST-GNN model across successive test chunks from various STKG datasets sourced from an open-source video dataset. The temporal coherence graph consistently demonstrates high stability, with coherence scores predominantly ranging from 0.94 to 0.95. This indicates that our model effectively preserves prediction consistency across sequential temporal inputs. It is evident that coherence experiences slightly more fluctuations for the Something-Something dataset, reflecting its intrinsic complexity due to the involvement of nuanced semantic actions and minimal temporal dynamics. In parallel, the windowed forgetting rate shows a significant decline over epochs, confirming the model’s learning stability as it effectively mitigates catastrophic forgetting across all datasets. Likewise, the Something-Something dataset displays a higher initial and sustained forgetting rate, yet it gradually improves by reinforcing our model’s strong incremental adaptation capability even in the face of challenging continual learning scenarios.

C. Model Implementation and Experimental Setup

All experiments were implemented in Python 3.8 with PyTorch 3.9 and PyTorch Geometric. We fixed the random seed to 42 and PyTorch to ensure reproducibility. We also share our custom STKGs as publicly available with augmented and a normalized timestamp feature file for real-time simulation. These timestamps are used to create time-aware chunks of size 300 with 10% overlap. These chunks are split chronologically into training (80%), validation (10%), and test (10%) sets.

Our CAST-GNN model uses two sequential Temporal Embedding layers (hidden dimension = 128, time embedding dimension = 32) followed by a 4-headed spatio-temporal TransformerConv. We incorporate a GRU-based EGPM and a TGNMemory module for node-level state retention, and an AdaptiveMetaLearning layer. Edge features are concatenated with time-difference embeddings as TransformerConv input.

Moreover, we used the AdamW optimizer with an initial learning rate of 1×10^{-3} and a weight decay of 1×10^{-5} , employed gradient clipping with a maximum norm of 0.8, and applied a StepLR scheduler to decay the learning rate by a factor of 0.1 every five epochs. Scheduler learning usage is only applied for the Something-Something dataset. For this reason, we share the model file of this dataset separately. Since the structures of the other datasets are the same, it is sufficient to share a common model file to process them by only adjusting the dataset path. Training was conducted for 10 epochs with meta-learning updates enabled (meta-learning

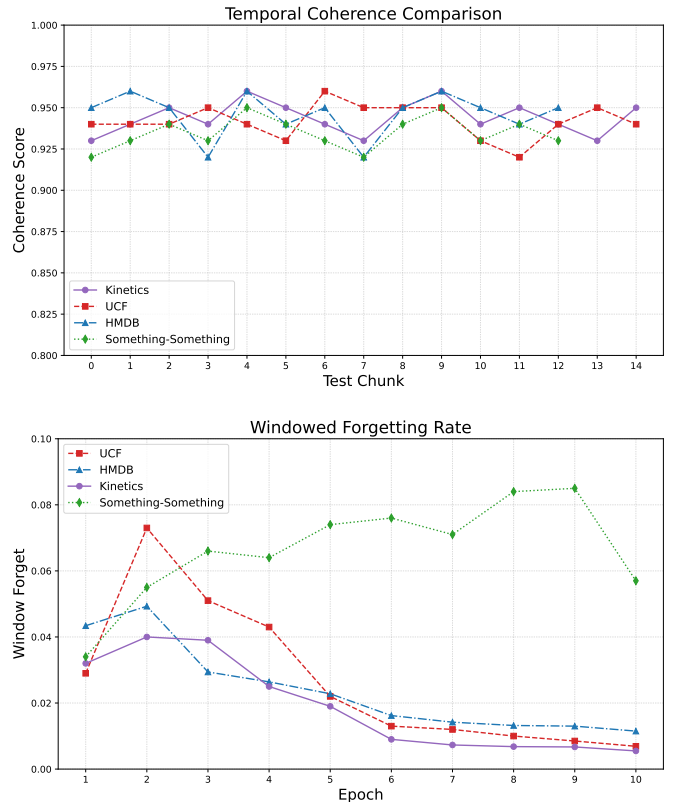


Fig. 2: Our proposed model’s temporal coherence trend against forgetting rate across chunks over various datasets which are used to derive our STKGs.

rate 1×10^{-3}), knowledge distillation ($\alpha = 0.3$, $T = 1.0$), and a Hybrid Selective Replay Buffer of capacity 300 using the “uncertainty-class” sampling strategy. Fisher information for Elastic Weight Consolidation was precomputed on the validation set prior to continual training.

All runtime measurements, both training and evaluation were executed on the CPU to be close real-time systems with compatible computational cost request on a machine with an 11th Gen Intel(R) Core(TM) i5-11300H @ 3.10GHz CPU.

D. Ablation Study

To comprehensively analyze the contribution and impact of the different components in our proposed CAST-GNN model, we conduct an ablation study. This analysis systematically examines how each specific module influences overall model performance in terms of F1 score, average forgetting rate, adaptation speed, and inference latency. We further investigate the effectiveness of the Hybrid Selective Replay Buffer by evaluating various sampling strategies to see how each strategy aligns with our continual learning goals.

The first ablation study is summarized in Table III, which demonstrates the critical roles of each module within our CAST-GNN architecture through our STKG data derived from Kinetics. We preferred this dataset for the ablation study because we thought that its complex structure would provide

TABLE III: Ablation study with the main components of our approach on the Kinetics Dataset

Ablation Module	(%) F1	(%) Average Forgetting Rate	Average Adaptation Speed (epochs)	Average Inference Latency (ms)
Full Method	97	0.17	8.00	10.06
(-) Self Attention Layer	70	0.06	8.76	8.24
(-) Meta Learning	83	0.20	8.63	10.83
(-) TGNMemory	40	1.55	7.63	10.35
(-) Knowledge Distillation	85	0.14	8.44	10.47
(-) Temporal Layer	76	0.07	8.83	9.47
(-) EGPM	47	0.42	7.62	3.75

TABLE IV: Ablation study with Hybrid Selective Sampling

Sampling Mode	% Accuracy	% F1	Recall@5	MRR
Random	81.5	80.68	99.62	89.52
Uncertainty	83.41	82.74	99.56	90.69
Class-balanced	81.23	80.43	98.26	88.60
<i>Uncertainty-class</i>	84.33	83.39	99.79	91.32
Proportional	79.33	77.79	99.18	88.07
Diversity	76.10	74.17	98.18	85.60

a better opportunity to observe the effects of the model components. The results show that the full model configuration achieves the highest performance with a 97% F1 score by highlighting the effectiveness of the integrated approach for robust continual learning. When we remove the Self-Attention Layer, the F1 score significantly drops to 70%. This indicates the importance of capturing spatio-temporal relationships between nodes by aligning with previous literature findings which emphasize on attention mechanisms for temporal dynamics [3]. On the other hand, removing the Meta-Learning component decreases the F1 score to 83% and slightly increases in forgetting rate. This showcases its effectiveness in quickly adapting to changes. On the other hand, removing one of the most important components, TGNMemory, drastically impacts performance. This clearly illustrates the essential role of maintaining temporal contextual memory to handle incremental updates effectively. Similarly, the EGPM proves its critical existence in performance stability by causing decrease in F1 to 47%. Knowledge Distillation and the Temporal Embedding Layer also notably influence the model performance, as their exclusion results in lower F1 scores. Also, these modules emphasize their role in mitigating the forgetting rate by enhancing semantic representations. Unlike other modules, the inference latency drops to 3.75 ms only in the absence of EGPM, reflecting the computational intensity but also the essential nature of this component for temporal consistency and robustness. Overall, these results affirm that each module provides exclusive capabilities by enabling our model to achieve superior accuracy, adapt quickly, and minimize catastrophic forgetting.

Table IV explores various selective sampling strategies within the Hybrid Selective Replay Buffer module in our method, to identify the best suited sampling mode for real-time continual learning scenarios. We conducted these ex-

periments with a small amount of STKG data derived from the Kinetics dataset by a quick training phase. Since our aim was only to observe the behavior of these modules, we interrupted the model training for each sampling mode and observed and reported its effect on the model via metrics. For this reason, this table does not show the maximum training performance of the model. We compare standard sampling methods: Random, Uncertainty-based [27], Class-balanced [28], Proportional [29], Diversity-based [30], and our proposed *Uncertainty-class* hybrid method. Random sampling serves as a baseline by achieving an accuracy of 81.5% and an F1 score of 80.68%. The uncertainty-based method, slightly improves accuracy (83.41%) and F1 (82.74%) by showing the benefit of retaining uncertain and informative samples. Class-balanced sampling achieves moderate results by demonstrating limited effectiveness in handling data distribution in streaming contexts. Similarly, the proportional strategy exhibits reduced performance by indicating its ineffectiveness in managing catastrophic forgetting. Lowest performance is obtained by Diversity-based sampling despite achieving high recall@5 (98.18%). This tells that diversity alone might neglect critical temporal and semantic patterns that are necessary for accurate classification. However, our proposed state-of-the-art *Uncertainty-class* strategy significantly outperforms all others, by reaching the highest accuracy (84.33%) and F1 score (83.39%), as well as superior recall@5 (99.79%) and MRR, 91.32%. This hybrid sampling approach optimally balances the importance of uncertainty-based informative samples and class distribution. Thus, it effectively enhances catastrophic forgetting by keeping high classification accuracy across incremental learning stages. These findings clearly indicate that the Hybrid Selective Replay Buffer with *Uncertainty-class* sampling strategy strongly complements of our CAST-GNN model, while other strategies either underperform due to neglecting class-balance or temporal coherence.

Collectively, these ablation studies demonstrate the importance and effectiveness of our carefully integrated CAST-GNN components and highlight the crucial role of optimized sampling strategies in enhancing real-time continual learning performance for spatio-temporal knowledge graphs.

VI. DISCUSSION

Our CAST-GNN model’s unified continual learning architecture demonstrably advances the state of the art in real-

time spatio-temporal reasoning. In the absence of any directly comparable continual-learning methods for streaming video-derived STKGs, we adapted and re-implemented the most related frameworks under equivalent conditions to provide a precise comparative analysis of our model’s performance. Across four diverse video-derived STKG benchmarks, our model consistently outperforms existing similar methods, achieving an average accuracy of 96–97% with forgetting rates between 0.13–31% and adaptation speeds under 9 epochs (Table II). The integration of temporal embedding layers and adaptive self-attention proves critical by ablating either component yields accuracy drops of over 15% and marked increases in forgetting (Table II). These results confirm that capturing both local temporal patterns and global structural context is essential for sustained performance in streaming environments.

However, CAST-GNN brings higher latency (10 ms/chunk vs. 3 ms for lighter methods), indicating a need for optimization through quantization or adaptive gating. Performance on fine-grained datasets like Something–Something remains moderate, so potential improvements via hierarchical attention or dynamic graph refinement should be explored. This performance degradation can be attributed to the dataset’s subtle temporal cues and fine-grained object interactions, which may not be fully captured by the current attention mechanism or memory modules. The limited context windows in such scenarios may cause the model to overlook critical semantic shifts.

Furthermore, our XAI-driven evaluation in Figure 1 confirms that CAST-GNN successfully aligns global and local feature attributions even when concept drift occurs. However, this interpretability pipeline currently depends on offline Integrated Gradient computations, which limits its applicability to real-time settings. Therefore, future research should investigate lightweight, incremental attribution methods that can operate online.

By demonstrating that integrating temporal embeddings with adaptive self-attention and hybrid replay strategies yields substantial improvements in continual semantic reasoning, CAST-GNN sets a new direction for unified architectures in graph-based continual learning. This holistic approach offers a foundation for future models that must balance adaptivity, knowledge retention, and explainability in real-world streaming applications. The proposed framework holds significant promise for practical deployment in areas such as real-time surveillance, environmental monitoring, and smart transportation systems, where semantic drift and non-stationarity data are prevalent challenges.

VII. CONCLUSION

In this work, we have presented CAST-GNN, the first comprehensive continual learning framework tailored for spatio-temporal knowledge graphs derived from real-time video streams. By unifying temporal embedding layers, adaptive self-attention, episodic graph pattern memory, hybrid selective replay buffers, and dual regularization (Fisher information and

knowledge distillation), CAST-GNN simultaneously addresses four core challenges: (i) incremental embedding adaptation, (ii) catastrophic forgetting mitigation, (iii) semantic-level reasoning, and (iv) interpretability.

This work establishes the first unified framework that bridges the gap between structural adaptation and semantic-level reasoning in STKGs, laying the groundwork for the next generation of explainable, online graph learning systems. Extensive experimentation on UCF101, HMDB51, Kinetics400, and Something–Something STKGs demonstrates its superior accuracy, low forgetting rates, rapid adaptation, and robust generalization across diverse domains ranging from activity recognition to urban mobility forecasting.

For future work, we aim to improve CAST-GNN’s performance and adaptability by implementing dynamic management of the memory buffer. Additionally, we plan to develop real-time explainability (XAI) methods that can instantly provide clear, node-level insights, enhancing the transparency of model decisions in globally important areas such as video analytics or environmental monitoring.

Given its capacity to operate under realistic latency constraints and its strong semantic reasoning abilities, CAST-GNN is well-positioned for integration into safety-critical systems where trust, transparency, and real-time responsiveness are paramount. Ultimately, we envision CAST-GNN serving as a foundation for autonomous agents capable of continuous understanding and reasoning over dynamic environments in an interpretable and reliable manner.

ACKNOWLEDGMENT

This work was supported by the Swiss National Science Foundation through the StreamKG project with grant number 213369.

REFERENCES

- [1] Y. Ma, Z. Guo, Z. Ren, J. Tang, and D. Yin, “Streaming graph neural networks,” in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, ser. SIGIR ’20. ACM, Jul. 2020, pp. 719–728.
- [2] E. Rossi, B. Chamberlain, F. Frasca, D. Eynard, F. Monti, and M. Bronstein, “Temporal graph networks for deep learning on dynamic graphs,” arXiv, 2020. [Online]. Available: <https://doi.org/10.48550/arXiv.2006.10637>
- [3] A. Sankar, Y. Wu, L. Gou, W. Zhang, and H. Yang, “Dysat: Deep neural representation learning on dynamic graphs via self-attention networks,” in *Proceedings of the 13th ACM International Conference on Web Search and Data Mining*. WSDM, 2020, pp. 519–527.
- [4] G. H. Nguyen, J. B. Lee, R. A. Rossi, N. K. Ahmed, E. Koh, and S. Kim, “Continuous-time dynamic network embeddings,” in *Companion of the The Web Conference 2018 on The Web Conference 2018 - WWW ’18*, ser. WWW ’18. ACM Press, 2018, pp. 969–976.
- [5] J. Wu, Y. Xu, Y. Zhang, C. Ma, M. Coates, and J. C. K. Cheung, “Tie: A framework for embedding-based incremental temporal knowledge graph completion,” in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, ser. SIGIR ’21. ACM, Jul. 2021, pp. 428–437.
- [6] M. Mirtaheeri, M. Rostami, and A. Galstyan, “History repeats: Overcoming catastrophic forgetting for event-centric temporal knowledge graph completion,” arXiv, pp. 7740–7755, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2305.18675>

- [7] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, Mar. 2017.
- [8] Z. Tian, D. Zhang, and H.-N. Dai, "Continual learning on graphs: A survey," arXiv, 2024. [Online]. Available: <https://doi.org/10.48550/arXiv.2402.06330>
- [9] Q. Yuan, S. Guan, P. Ni, T. Luo, K. Man, P. W. H. Wong, and V. I. Chang, "Continual graph learning: A survey," arXiv, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2301.12230>
- [10] R. Wang, J. Huang, Y. Zhang, J. Li, Y. Wang, W. Zhao, S. Liu, C. Mendis, and T. F. Abdelzaher, "Tgonline: Enhancing temporal graph learning with adaptive online meta-learning," in *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 7 2024, pp. 1659–1669.
- [11] C. Wang, G. Tan, S. B. Roy, and B. Ooi, "Distribution-aware online continual learning for urban spatio-temporal forecasting," 2024. [Online]. Available: <https://doi.org/10.48550/arXiv.2411.15893>
- [12] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre, "Hmdb: A large video database for human motion recognition," in *2011 International Conference on Computer Vision*. IEEE, Nov. 2011, pp. 2556–2563.
- [13] K. Soomro, A. Zamir, and M. Shah, "Ucf101: A dataset of 101 human actions classes from videos in the wild," arXiv, 2012. [Online]. Available: <https://doi.org/10.48550/arXiv.1212.0402>
- [14] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman, "The kinetics human action video dataset," arXiv, 2017. [Online]. Available: <https://arxiv.org/abs/1705.06950>
- [15] R. Goyal, S. E. Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag, F. Hoppe, C. Thureau, I. Bax, and R. Memisevic, "The "something something" video database for learning and evaluating visual common sense," in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2017.
- [16] G. Tataroğlu Özbülak, Y. Shrestha, and J.-P. Calbimonte, "Stkgnn: Scalable spatio-temporal knowledge graph reasoning for activity recognition," in *Proceedings of the 34th ACM International Conference on Information and Knowledge Management CIKM '25 (in press)*, 2025.
- [17] S. Kumar, X. Zhang, and J. Leskovec, "Predicting dynamic embedding trajectory in temporal interaction networks," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 1269–1278.
- [18] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," in *Advances in Neural Information Processing Systems*, vol. 30, 2017, pp. 6467–6476.
- [19] G. E. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," arXiv, 2015. [Online]. Available: <https://arxiv.org/abs/1503.02531>
- [20] R. Tiwari, K. Killamsetty, R. Iyer, and P. Shenoy, "Gcr: Gradient coreset based replay buffer selection for continual learning," in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jun. 2022, pp. 99–108.
- [21] P. Arora, Deepali, and S. Varshney, "Analysis of k-means and k-medoids algorithm for big data," *Procedia Computer Science*, vol. 78, pp. 507–512, 2016.
- [22] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2017.
- [23] A. Pareja, G. Domeniconi, J. Chen, T. Ma, T. Suzumura, H. Kanezashi, T. Kaler, T. Schardl, and C. Leiserson, "Evolvegcn: Evolving graph convolutional networks for dynamic graphs," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, pp. 5363–5370, Apr. 2020.
- [24] A. Chaudhry, P. Dokania, T. Ajanthan, and P. H. S. Torr, "Riemannian walk for incremental learning: Understanding forgetting and intransigence," in *Proceedings of the European conference on computer vision (ECCV)*. Springer International Publishing, 2018, pp. 556–572.
- [25] J. Gama, I. Žliobaitė, A. Bifet, M. Pechenizkiy, and A. Bouchachia, "A survey on concept drift adaptation," *ACM Computing Surveys*, vol. 46, no. 4, pp. 1–37, Mar. 2014.
- [26] S. Akkas and A. Azad, "Gnnshap: Scalable and accurate gnn explanation using shapley values," in *Proceedings of the ACM Web Conference 2024*, 2024.
- [27] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *Proceedings of The 33rd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. F. Balcan and K. Q. Weinberger, Eds., vol. 48. New York, New York, USA: PMLR, 20–22 Jun 2016, pp. 1050–1059.
- [28] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jul. 2017, pp. 5533–5542.
- [29] F. M. Castro, M. J. Marín-Jiménez, N. Guil, C. Schmid, and K. Alahari, "End-to-end incremental learning," in *Proceedings of the European conference on computer vision (ECCV)*. Springer International Publishing, 2018, pp. 241–257.
- [30] R. Aljundi, M. Lin, B. Goujaud, and Y. Bengio, "Gradient based sample selection for online continual learning," in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. Red Hook, NY, USA: Curran Associates Inc., 2019.