

Differentiable Cluster Discovery in Temporal Graphs

Md Hafizur Rahman

*Department of Mechanical and Industrial Engineering
University of Toronto*

hafizur.rahman@mail.utoronto.ca

Chi-Guhn Lee

*Department of Mechanical and Industrial Engineering
University of Toronto*

chiguhn.lee@utoronto.ca

Reviewed on OpenReview: <https://openreview.net/forum?id=1caZVb6zL7>

Abstract

Existing temporal graph clustering methods suffer from poor optimization dynamics due to reliance on heuristically initialized cluster assignment distribution without considering the dynamic nature of the evolving graph. The target cluster assignment distribution often conflicts with evolving temporal representations, leading to oscillatory gradients and unstable convergence. Motivated by the need for differentiable and adaptive clustering in dynamic settings, we propose TGRAIL (**T**emporal **G**raph **A**lignment and **I**ndex **L**earning), a novel end-to-end framework for temporal graph clustering based on Gumbel–Softmax sampling. TGRAIL enables discrete cluster assignments while maintaining the gradient flow. To ensure stable training, we formulate the clustering objective as an expectation over Monte Carlo samples and show that this estimator is both unbiased and variance-reduced. Furthermore, we incorporate a temporal consistency loss to preserve the order of interactions across time. Extensive experiments on six real-world temporal graph datasets demonstrate that our approach consistently outperforms state-of-the-art baselines, achieving higher clustering accuracy and robustness. Our results validate the effectiveness of jointly optimizing temporal dynamics and discrete cluster assignments in evolving graphs.

1 Introduction

Graphs are fundamental tools for modeling relationships and interactions in complex systems, spanning domains such as social networks, biological networks, communication systems, and financial markets (Hamilton et al., 2017; Ying et al., 2019; Sun et al., 2020; Wang et al., 2022; Zheng et al., 2025). A central task in graph analysis is clustering, which aims to group nodes into communities based on structural or semantic similarity. Traditional graph clustering methods operate on static graphs, where the topology and node attributes remain fixed. These methods, including spectral clustering and modularity-based approaches (Bianchi et al., 2020; You et al., 2021; Tsitsulin et al., 2023; Liu et al., 2025a), have been widely adopted due to their theoretical foundations and interpretability.

Compared to this, deep clustering methods integrate representation learning with clustering objectives. For instance, Deep Embedded Clustering (DEC) (Xie et al., 2016) combines autoencoder-based embeddings with Kullback–Leibler (KL) divergence-based soft assignments. Extensions such as Improved DEC (Guo et al., 2017) and Structural Deep Clustering Networks (SDCN) (Bo et al., 2020) incorporate reconstruction losses or graph neural networks to better leverage node features and topology. Despite their success, these methods are fundamentally static: they assume access to a complete adjacency matrix and cannot model temporal dependencies. Consequently, they are unable to capture the evolving nature of communities or adapt to dynamic patterns of interaction.

Recently, temporal graph clustering has emerged and gained attention to address these limitations. A temporal graph captures the temporal dimension through a sequence of time-stamped events. Instead of

modeling edges as static relations, temporal graphs represent interactions as sequences, allowing finer-grained analysis of how relationships form, persist, and dissolve over time. This richer representation enables new opportunities, such as tracking evolving communities, detecting temporal anomalies, and forecasting future events (Cini et al., 2025; Postuvan et al., 2024; Liu et al., 2024).

Several approaches have been proposed to model temporal graphs. Time-aware graph neural networks (TGNNs), such as TGAT (Xu et al., 2020), TGN (Rossi et al., 2020), and HTNE (Zuo et al., 2018), introduce temporal attention, memory, or Hawkes processes to encode evolving features. However, these methods typically decouple representation learning from clustering, requiring a post hoc clustering step. This two-stage design can be suboptimal, as the learned representations may not align well with the clustering objective, and errors from the first stage propagate without correction. Moreover, the clustering step is non-differentiable, preventing end-to-end training.

Recent methods attempt to address this limitation by integrating clustering within the training loop. For example, TGC (Liu et al., 2024) incorporates a clustering loss into the temporal graph encoder using soft assignments derived from a Student’s t -distribution. This approach enables joint optimization of embeddings and cluster centroids. While the target distribution is expressed as time dependent in their approach, its reliance on fixed node embeddings results in a distribution that does not evolve over time which fails to adapt temporally consistent cluster assignment. Additionally, the t -distribution has several drawbacks in dynamic settings: it assumes a fixed degree-of-freedom parameter, is sensitive to initialization, and tends to overemphasize outliers due to its heavy-tailed nature (Linderman & Steinerberger, 2019). Figure 1 we provide a visual depiction of temporal cluster dynamics in evolving graphs.

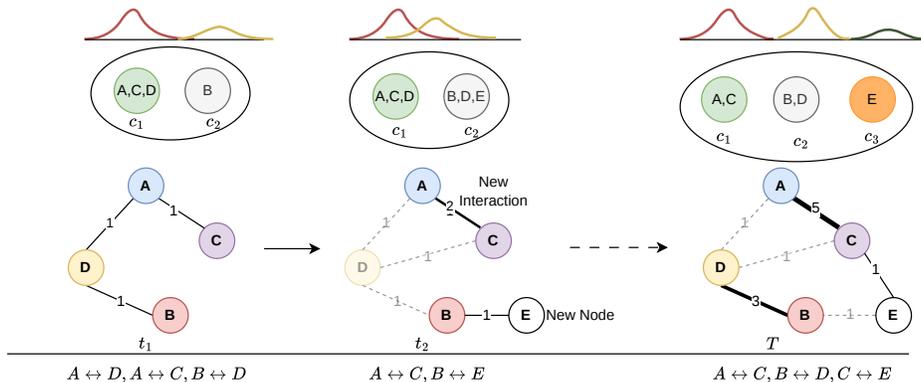


Figure 1: Temporal evolution of cluster assignments in dynamic graphs. Nodes may shift clusters due to new interactions, inactive nodes, or structural changes over time.. At t_1 , nodes A, C, D form cluster C_1 , and B belongs to C_2 . A new node E appears at t_2 , reshaping interactions and leading to reassignment of B, D, E to C_2 . By time T , node E forms a separate cluster C_3 .

To address the aforementioned limitations, we propose a novel, differentiable framework for temporal graph clustering. We formulate the cluster assignment process as stochastic sampling from a Gumbel-Softmax distribution, which enables discrete assignments to be learned through gradient-based optimization. We summarize our contributions as follows-

1. **A differentiable framework for temporal graph clustering.** We propose TGRAIL, a method that *jointly* learns node representations and discrete cluster assignments in dynamic graphs via a Monte Carlo Gumbel Softmax re-parameterization. This removes the need for *post-hoc* process or t -distribution soft assignments, enabling end-to-end training thus aligns cluster assignment with temporal node embeddings.
2. **Unbiased, low-variance gradient estimation with theoretical guarantees.** We derive a tight variance bound for the Gumbel estimator and prove a non-asymptotic SGD convergence theorem under

standard Lipschitz and bounded-step assumptions. Our analysis clarifies why discrete assignments remain stable throughout training.

3. **A unified temporal–clustering loss that scales linearly in interactions.** By coupling a temporal contrastive objective with the discrete clustering loss, we keep complexity at $\mathcal{O}(|E|)$ rather than $\mathcal{O}(N^2)$, making TGRAIL practical for long, sparse interaction streams.
4. **Extensive empirical validation on six evolving-graph datasets.** TGRAIL outperforms ten SOTA baselines by 3–5% macro- F_1 on sparse datasets (PATENT, DBLP) and matches or exceeds the best methods on dense or highly non-stationary graphs.

2 Related Work

Graph Clustering via Neural Networks. Initial deep learning approaches to graph clustering leveraged MLP-based autoencoders to extract latent node embeddings from the graph structure. GraphEncoder (Tian et al., 2014) and DNCR (Cao et al., 2016) encoded proximity between nodes using sparse autoencoders and random walk-based techniques, followed by k -means clustering. These early methods demonstrated that deep representations could improve clustering performance but struggled to integrate node attribute information. The introduction of graph convolutional networks enabled models to jointly encode structural and attribute information. Kipf and Welling’s VGAE (Kipf & Welling, 2016) and Wang et al.’s MGAE (Wang et al., 2017) used graph encoders to produce informative latent spaces for downstream clustering. These works laid the groundwork for reconstructive methods (Wang et al., 2019; Park et al., 2019; Bandara et al., 2024; Yu et al., 2025) where reconstruction of adjacency or feature matrices acted as the self-supervised objective. Similarly, adversarial mechanisms were introduced to regularize latent spaces and improve representation robustness. ARGAN (Pan et al., 2018) employed a discriminator to align latent embeddings with a Gaussian prior, while CommunityGAN (Jia et al., 2019) generated synthetic samples for structure-preserving embedding. Though effective in reducing overfitting and capturing community semantics, these methods separate clustering from representation learning, which introduces unstable cluster assignments.

Clustering-Oriented Architectures and Fusion Models. On the other hand, several methods sought to unify representation learning with clustering objectives. DAEGC (Wang et al., 2019) proposed attention-based graph encoders guided by clustering alignment loss. GALA (Park et al., 2019) enhanced encoder expressiveness via Laplacian sharpening. Models like SDCN (Bo et al., 2020) and DFCN (Tu et al., 2021) integrated attribute and structure views using novel fusion strategies, demonstrating that explicit clustering supervision during representation learning improved cluster separation and compactness. As graph sizes increased, scalability became a central concern (Xu et al., 2025). S3GC (Devvrit et al., 2022) performed scalable contrastive learning using batch-wise subgraph sampling and post-hoc k -means clustering. Dink-Net (Liu et al., 2023) unified contrastive representation learning and clustering optimization via differentiable dilation and shrinkage losses, enabling end-to-end training on graphs with over 100M nodes.

Dynamic/Temporal Graph Clustering. Temporal graph clustering extends conventional graph clustering to dynamic scenarios where node interactions evolve over time. Liu et al. (Liu et al., 2024) propose a general framework called Temporal Graph Clustering (TGC). This framework integrates temporal representation learning with clustering objectives tailored for interaction-sequence data. MVTGC (Liu et al., 2025b), CGC (Park et al., 2022) utilize multiview and contrastive objectives between graph snapshots to capture evolving community structures. These models address the temporal nature of clustering, which static methods cannot handle effectively.

Despite these advances, GNN-based temporal graph clustering approaches model t -distribution as a cluster assignment distribution (Bo et al., 2020; Liu et al., 2024; 2025a), which may be suboptimal in dynamic settings due to its heavy tails that amplify the influence of transient or noisy nodes. In contrast to these approaches, we learn the cluster assignment using Gumbel Softmax in an end-to-end manner, aligning temporal evolution and cluster assignment.

3 Temporal Graph Clustering

3.1 Problem Definition

As stated in the previous section, temporal graphs capture not a fixed structure but an evolving stream of interactions. In such dynamic networks, whether social platforms, citation graphs, or sensor grids where nodes can emerge, disappear, or reconfigure their connections over time. This evolution manifests as fluctuations in node activity, shifting neighborhood contexts, and changing roles, all of which influence the cluster membership of each node at every timestamp. To capture this temporal evolution, we consider the network as a continuous sequence of graphs where the topology is a chronological stream of temporal events Rossi et al. (2020); Chmura et al. (2025). Let the sequence be $\{G^{(1)}, G^{(2)}, \dots, G^{(T)}\}$, where each timestamp $G^{(t)} = (\mathcal{V}^{(t)}, \mathcal{E}^{(t)})$ represents the network’s state at time t where $\mathcal{V}^{(t)}$ denotes the set of active nodes, and $\mathcal{E}^{(t)} \subseteq \mathcal{V}^{(t)} \times \mathcal{V}^{(t)}$ defines their pairwise interactions. We can define the problem of temporal graph clustering as follows. For notation clarity, we denote matrices in bold capital letters, vectors in bold small letters, and scalars in non-bold letters.

Problem 3.1 (Continuous-Time Temporal Graph Clustering). Given a temporal graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{T})$ and time-dependent node features $\mathbf{X}^{(t)} \in \mathbb{R}^{N \times D}$ and adjacency matrix $\mathbf{A}^{(t)} \in \mathbb{R}^{N \times N}$ at each timestamp $t \in \mathcal{T}$, the objective is to learn a continuous-time node encoder f_θ and cluster centroids $\mathbf{C}^{(t)} = \{\mathbf{c}_i^{(t)} \dots \mathbf{c}_K^{(t)}\}$ parameterized by an assignment mechanism q_ϕ , such that the learned soft assignments exhibit compact cluster structure. Specifically, given the historical interaction upto time t , we aim to learn,

$$\mathbf{Z}^{(t)} = f_\theta(\mathbf{X}^{(\leq t)}, \mathbf{A}^{(\leq t)}) \quad ; \quad \mathbf{\Pi}^{(t)} = q_\phi(\mathbf{Z}^{(t)}). \quad (1)$$

Here, $\mathbf{Z}^{(t)}$ is the latent embedding matrix, and $\mathbf{\Pi}^{(t)} = [\boldsymbol{\pi}_1^{(t)}, \dots, \boldsymbol{\pi}_N^{(t)}]$ is the cluster assignment matrix, where each $\boldsymbol{\pi}_i^{(t)}$ is a soft cluster membership vector for node i at time t , lying on the $(K-1)$ -dimensional probability simplex defined as-

$$\Delta^{K-1} := \left\{ \boldsymbol{\pi}_i^{(t)} \in \mathbb{R}^K \left| \sum_{k=1}^K \pi_{i,k} = 1 \text{ and } \pi_{i,k} \geq 0 \text{ for all } k \right. \right\}. \quad (2)$$

3.2 Joint Representation Learning and Clustering Objective

Building on our temporal graph formulation, from Equation 1, it is evident that the temporal graph clustering problem naturally lends itself to an optimization formulation, where we need to simultaneously optimize node representations and cluster assignments while maintaining temporal consistency. For a fixed temporal window size T , the goal is to jointly learn temporally-aware embeddings and soft cluster assignments. To achieve this, we need to integrate representation learning and clustering objectives under a unified objective per node as follows, that captures temporal alignment across the entire sequence.

$$\min_{\theta, \phi} \sum_{t=1}^T \mathbb{E}_{\mathbf{x}_i^{(t)} \sim p_{\text{data}}(\mathbf{x}_i^{(t)})} \left[\mathbb{E}_{\boldsymbol{\pi}_i^{(t)} \sim q_\phi(\cdot | \mathbf{z}_i^{(t)})} \mathcal{L}_{\text{clu}}(\mathbf{x}_i^{(t)}, \mathbf{z}_i^{(t)}, \boldsymbol{\pi}_i^{(t)}) \right] \quad (3)$$

Here, \mathcal{L}_{clu} is a clustering loss function that evaluates the quality of the assignments $\boldsymbol{\pi}_i^{(t)}$ based on the latent embeddings and their temporal consistency. The outer expectation captures variability in the input, while the inner expectation reflects the stochasticity of cluster assignments. Bo et al. (2020); van der Maaten & Hinton (2008); Liu et al. (2024) employs soft clustering methods using the Student- t distribution to define the cluster assignment probability vector $\boldsymbol{\pi}_i$ for a node, especially in deep embedding-based approaches. Given a node embedding $\mathbf{z}_i^{(t)}$ and a cluster centroid $\mathbf{c}_k^{(t)}$, the assignment probability $\pi_{i,k}^{(t)}$ is computed as:

$$\pi_{i,k}^{(t)} = \frac{(1 + \|\mathbf{z}_i^{(t)} - \mathbf{c}_k^{(t)}\|^2 / \nu)^{-\frac{\nu+1}{2}}}{\sum_{j=1}^K (1 + \|\mathbf{z}_i^{(t)} - \mathbf{c}_j^{(t)}\|^2 / \nu)^{-\frac{\nu+1}{2}}} \quad (4)$$

Here, ν is the degrees of freedom (commonly set to 1), and the distribution emphasizes local structure by assigning higher probability to closer centroids while retaining robustness to outliers due to its heavy-tailed nature. To improve convergence and increase assignment confidence, a sharpened target distribution (Bo et al., 2020; Liu et al., 2024) $\tilde{\pi} = \{\tilde{\pi}_{i,1} \dots \tilde{\pi}_{i,K}\}$ is computed by squaring and normalizing the initial assignments, and the following is defined as clustering loss as Kullback–Leibler (KL) divergence to jointly update the node embeddings and centroids.

$$\mathcal{L}(\theta, \phi) = KL(\pi_i^{(t)} || \tilde{\pi}) \quad (5)$$

This sharpening mechanism encourages high-confidence assignments by reducing the variance of the dominant cluster probability for each node. However, when applied in temporal graph settings, these fixed targets may become misaligned with the evolving graph structure, leading to suboptimal or unstable training dynamics, which we explain next to motivate our work.

3.3 Challenges: Gradient Conflicts in Temporal Clustering

Optimizing the clustering objective in Equation 5 involves updating both the encoder parameters θ and the centroid centroids, where the loss is defined as the KL divergence between the current assignment $\pi_{i,k}^{(t)}$ and the sharpened target $\tilde{\pi}_{i,k}$. Taking the gradient of the KL loss with respect to the node embedding induces a force (derivation is given in the Appendix 9):

$$F_{i,k}^{(t)} = \underbrace{\frac{2\pi_{i,k}^{(t)}d_{i,k}^{(t)}}{1 + (d_{i,k}^{(t)})^2}}_{\text{Geometric term } G(d,\pi)} \cdot \left[\underbrace{(\pi_{i,k}^{(t)} - \tilde{\pi}_{i,k})}_{\text{Target error}} + \underbrace{(\pi_{i,k}^{(t)} - 1) \log \pi_{i,k}^{(t)}}_{\text{Entropy regularization}} \right] \quad (6)$$

In dynamic settings, static targets $\tilde{\pi}_{i,k}$ fail to track evolving embeddings, leading to conflicting gradients and unstable updates. Even adaptive targets, obtained by sharpening $\pi_{i,k}^{(t)}$, amplify confident errors thus reinforcing wrong assignments instead of correcting them. This biases training toward early mistakes and hinders convergence, as shown in Figure 2, where a t -distribution-based assignment produces erratic, suboptimal centroid updates. In contrast, our proposed mechanism better aligns assignments with evolving communities which we describe below.

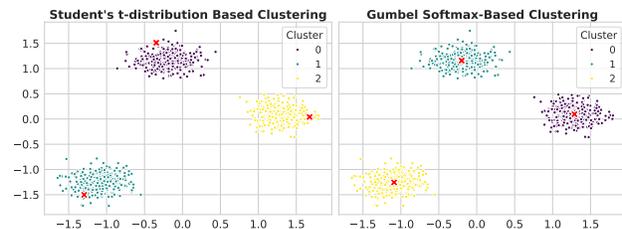


Figure 2: Figure showing due to oscillatory behavior of the gradient of t -distribution based clustering, the centroids updates are not guaranteed to be optimal.

4 Proposed Method

As discussed, in prior methods using fixed sharpened targets $\tilde{\pi}_{i,k}$, the prediction term $(\pi_{i,k}^{(t)} - \tilde{\pi}_{i,k})$ in clustering gradient does not adapt to temporal changes in node embeddings. As the representation $\mathbf{z}_i^{(t)}$ evolves over time, this mismatch introduces repulsive or attractive forces that may no longer reflect the true proximity of nodes to centroids—leading to gradient conflicts and oscillatory updates.

Therefore, we aim to remove this non-adaptive target by directly sampling the cluster assignment from the cluster assignment distribution and align the assignment according to the updated temporal node embeddings. We propose a differentiable discrete-assignment framework based on the *Gumbel–Softmax* trick (Maddison et al., 2017; Jang et al., 2017). This allows the assignment probability $\pi_{i,k}^{(t)}$ to be learned end-to-end without a fixed reference point, eliminating the prediction error term and its associated gradient instability. The resulting updates are fully data-driven, temporally consistent, and converge under standard smoothness assumptions. For every node $i \in \mathcal{V}^{(t)}$ we maintain a soft cluster-membership vector $\pi_i^{(t)} \in \Delta^{K-1}$ with entries $\pi_{i,k}^{(t)}$ (the probability that node i belongs to cluster k). Given unnormalised logits $\ell_{i,k}^{(t)} \in \mathbb{R}$ and i.i.d. noise

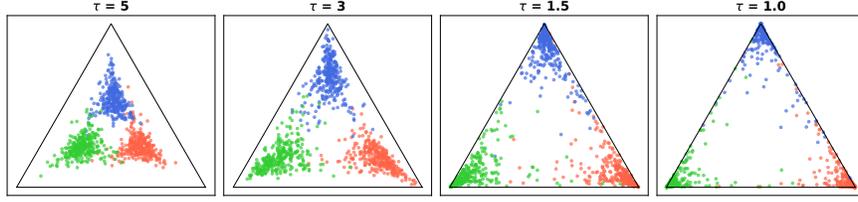


Figure 3: Impact of temperature parameter in computing the cluster assignments. For higher τ , soft assignments are smoother and more uniform across clusters which encourages exploration and better gradient flow, which is beneficial during early training when representations are still being learned. For smaller τ , the assignments become more discrete (closer to one-hot vectors), aligning better with the intended clustering objective.

variables $g_{i,k} \sim \text{Gumbel}(0, 1)$, the assignment distribution can be expressed as,

$$\pi_{i,k}^{(t)} = \frac{\exp((\log \ell_{i,k}^{(t)} + g_{i,k})/\tau)}{\sum_{j=1}^K \exp((\log \ell_{i,j}^{(t)} + g_{i,j})/\tau)}, \quad \tau > 0, \quad (7)$$

where the temperature τ controls discreteness as shown in Figure 3 ($\tau \rightarrow 0$ recovers hard one-hot vectors as the distribution becomes discrete). Given a collection of independent Gumbel noise variables \mathbf{g} , we can define soft cluster assignment as,

$$\mathbf{\Pi}^{(t)} = h_{\phi}(\mathbf{g}), \quad \text{where } \mathbf{g} \sim \text{Gumbel}(0, 1), \quad (8)$$

and $h_{\phi}(\cdot)$ is the Gumbel–Softmax mapping parameterized by ϕ . Given node embeddings $\mathbf{Z}^{(t)} = f_{\theta}(\mathbf{X}^{(t)})$ produced by the encoder f_{θ} and cluster centroids $\mathbf{C}^{(t)}$, we define the clustering objective as the expectation over the random Gumbel noise:

$$\mathcal{L}(\mathbf{X}^{(t)}, \mathbf{C}^{(t)}; \theta, \phi) := \mathbb{E}_{\mathbf{g} \sim \text{Gumbel}(0,1)} \left[\mathcal{L}_{\text{clu}}(f_{\theta}(\mathbf{X}^{(t)}), \mathbf{C}^{(t)}, h_{\phi}(\mathbf{g})) \right], \quad (9)$$

Since the expectation in Equation 9 involves nonlinear transformations of stochastic samples through the Gumbel–Softmax reparameterization $h_{\phi}(\mathbf{g})$ and the clustering loss \mathcal{L}_{clu} becomes intractable to compute in closed form. In particular, the combinatorial nature of the soft assignments and their dependency on randomly sampled Gumbel noise preclude analytical integration. Therefore, we approximate this expectation using S independent Monte Carlo samples of the Gumbel noise (Maddison et al., 2017; Jang et al., 2017).

$$\mathcal{L}(\mathbf{X}^{(t)}, \mathbf{C}^{(t)}; \theta, \phi) := \mathbb{E}_{\mathbf{g} \sim \text{Gumbel}(0,1)} \left[\mathcal{L}_{\text{clu}}(f_{\theta}(\mathbf{X}^{(t)}), \mathbf{C}^{(t)}, h_{\phi}(\mathbf{g})) \right] \quad (10)$$

$$= \mathbb{E}_{\mathbf{g}} \left[\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \pi_{i,k}^{(t)}(\mathbf{g}) \times d_{i,k}^{(t)} \right] \quad (11)$$

$$\approx \frac{1}{S} \sum_{s=1}^S \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \frac{\exp((\log \ell_{i,k}^{(t)} + g_{i,k}^{(s)})/\tau) \times d_{i,k}^{(t)}}{\sum_{j=1}^K \exp((\log \ell_{i,j}^{(t)} + g_{i,j}^{(s)})/\tau)} \quad g_{i,k}^{(s)} \stackrel{\text{i.i.d.}}{\sim} \text{Gumbel}(0, 1). \quad (12)$$

where $d_{i,k}^{(t)}$ is the distance between cluster k and node i at time t . Equation 10 demonstrates that the clustering loss can be approximated by drawing S independent Gumbel–Softmax samples, evaluating the loss for each sample, and averaging the results. Since Gumbel noise makes Equation 12 differentiable, it integrates seamlessly with backpropagation as both the encoder f_{θ} and cluster centroids receive gradients as if the assignments were continuous. A full algorithm to compute the clustering loss is given in Algorithm 1 in the Appendix 11.

Temporal-consistency loss. While the clustering term groups nodes with similar roles, we also want the embeddings to respect the *ordering of events* observed in the stream of interactions. We treat the similarity

between node embeddings as a proxy conditional intensity of an interaction. To achieve this, we use the embedding-similarity score to estimate the Hawkes intensity score (Zuo et al., 2018). Let $\mathcal{E}^{(t)} \subseteq \mathcal{V}^{(t)} \times \mathcal{V}^{(t)}$ be the set of observed interactions at time t , then we can measure the intensity between node i and j at time t as-

$$s(\mathbf{z}_i^{(t)}, \mathbf{z}_j^{(t)}) = s_\mu(\mathbf{z}_i^{(t)}, \mathbf{z}_j^{(t)}) + \sum_{\substack{h \in \mathcal{H}_i \\ t' < t}} \alpha_{hj} s_\alpha(\mathbf{z}_h^{t'}, \mathbf{z}_j^{(t)}) e^{-\delta_{hj}(t-t')} \quad (13)$$

Here, $s_\mu(\cdot)$ is the cosine similarity score of node i and target node j at current time t and $s_\alpha(\cdot)$ is the cosine similarity between target node j and source node i 's historical neighbors. α_{hj} is the importance weight, and δ_{hj} is the exponential decay parameter. The exponential decay smoothly diminishes the influence of historical neighbor interactions as they become more temporally distant. For timestamp t , we distinguish positive intensities for observed edges $(i, j) \in \mathcal{E}_t$ and negative intensities for non-interacting pairs (i, b) drawn by negative sampling and define the contrastive temporal loss as negative log-likelihood with B negative samples per positive pair-

$$\mathcal{L}_{\text{tem}}(\theta) = -\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{(i,j) \in \mathcal{E}^{(t)}} \left[\log \sigma(s(\mathbf{z}_i^{(t)}, \mathbf{z}_j^{(t)})) + \sum_{b=1}^B \log(1 - \sigma(s(\mathbf{z}_i^{(t)}, \mathbf{z}_{n_b}^{(t)}))) \right], \quad (14)$$

where $\{n_b\}_{b=1}^B$ are negative samples and $\sigma(\cdot)$ denotes the sigmoid function and $s(\cdot)$ is computed using Equation 13. Combining Equation 10 and 14, we get the overall objective function as,

$$\mathcal{J}(\theta, \phi) = \mathcal{L}(\mathbf{X}^{(t)}, \mathbf{C}^{(t)}; \theta, \phi) + \lambda \mathcal{L}_{\text{tem}}(\theta) \quad (15)$$

where $\lambda > 0$ trades off cluster compactness against temporal predictability. The clustering term \mathcal{L} organises the latent space into coherent communities, while the temporal term \mathcal{L}_{tem} keeps consecutive embeddings faithful to the observed interaction sequence.

Theoretical Analysis We establish that the Gumbel-Softmax optimization procedure employed in our framework converges to a stationary point of the temporal clustering objective under standard assumptions of smoothness and bounded variance. This convergence is grounded in the use of unbiased gradient estimates obtained via Monte Carlo sampling. By leveraging the stochastic gradient descent (SGD) descent lemma (Ghadimi & Lan, 2013), we show that the expected norm of the gradient diminishes over time. Furthermore, our analysis incorporates the annealing of the temperature parameter τ , which progressively sharpens the cluster assignments from soft to nearly discrete. A key result underpinning this behavior is Lemma 4.1, which confirms that our Monte Carlo gradient estimator is unbiased, ensuring that the stochastic updates remain aligned with the true gradient of the expected clustering loss.

Lemma 4.1 (Unbiasedness). Let $\hat{\mathbf{g}}$ be the Monte-Carlo gradient estimator in Equation 10; then $\mathbb{E}[\hat{\mathbf{g}}] = \nabla_{\Theta} \mathcal{L}(\Theta)$, where $\Theta = \{\theta, \phi\}$ denotes the collection of encoder and assignment parameters.

Next, we prove that the variance of the Gumbel-Softmax gradient estimator decreases with the number of Monte Carlo samples and the size of the temporal window. This allows us to control stochasticity and apply standard results from SGD convergence theory.

Theorem 4.1 (Variance Bound). If $\|\nabla_{\Theta} \mathcal{L}_{\text{clu}}(\mathbf{Z}^{(t)}, \mathbf{C}^{(t)}, \mathbf{\Pi}^{(t)})\| \leq G_{\text{max}}$ for all admissible $(\mathbf{Z}^{(t)}, \mathbf{C}^{(t)}, \mathbf{\Pi}^{(t)})$, then $\text{Var}[\hat{\mathbf{g}}] \leq \frac{G_{\text{max}}^2}{ST}$, where T is the temporal context length in each mini-batch.

Combining these results, we show that the expected gradient norm of the clustering objective vanishes over time (Theorem 4.1). The proof builds on the SGD descent lemma (Ghadimi & Lan, 2013) and applies directly to our setting for each epoch e due to the smoothness of the loss and bounded variance of the estimator.

Theorem 4.2 (Convergence of Gumbel-Softmax Assignment in Temporal Clustering). Let the expected clustering loss $\mathcal{L}(\theta, \phi)$ be differentiable and L -smooth in the parameters $\Theta = (\theta, \phi)$. Assume the Monte-Carlo gradient used in training is an unbiased estimator of $\nabla \mathcal{L}$ with bounded second moment. If stochastic gradient descent is run with a constant stepsize $\eta \leq 1/L$ (or any diminishing stepsize satisfying $\sum_e \eta_e = \infty$ and $\sum_e \eta_e^2 < \infty$), then the parameter sequence $\{(\theta^{(e)}, \phi^{(e)})\}_{e=1}^{\infty}$ generated by the algorithm obeys,

$$\lim_{E \rightarrow \infty} \frac{1}{E} \sum_{e=1}^E \mathbb{E} \left[\|\nabla \mathcal{L}(\Theta^{(e)})\|^2 \right] = 0.$$

By extending this theorem, we can show that in the annealed setting where temperature $\tau^e \rightarrow 0$ guides the model from soft assignments to discrete ones as the loss function remains smooth and its variance is bounded, which justifies our choice of exponential decay of the temperature parameter. The complete proofs are provided in the Appendix 10.

Complexity of Temporal Graph Clustering with Gumbel-Softmax. In a temporal setting, a feasible model must update itself on the fly without ever materialising the full $N \times N$ adjacency matrix. Any procedure whose cost scales as $\mathcal{O}(N^2)$ quickly becomes intractable, whereas an $\mathcal{O}(|\mathcal{E}|)$ routine can process the stream event-by-event and train in mini-batches on commodity hardware and the optimiser visits every interaction once, yielding $\mathcal{O}(|\mathcal{E}|)$ time and memory (Liu et al., 2024). Introducing Gumbel-Softmax leaves this asymptotic bound unchanged. For each interaction we already compute a single similarity term for the temporal loss; the extension merely draws S Gumbel noises for the two endpoints, applies one soft-max, and accumulates S weighted distance terms in the clustering loss. These additions are $\mathcal{O}(S)$ per event, and S is a small, fixed constant. In practice as increasing number of samples does not always guarantee better performance (Paulus et al., 2021; Rainforth et al., 2018). Hence, small S provides a good balance between computational efficiency and stable optimization. Now, aggregated over the full sequence, the runtime becomes $c_1|\mathcal{E}| + c_2K|\mathcal{E}| = \mathcal{O}(|\mathcal{E}|)$. Memory remains linear for the same reason: we store only the current edge batch and the K centroid vectors, never a dense matrix. Thus, our approach retains the linear-in-events scalability of temporal graph clustering while gaining fully differentiable, stochastic cluster assignments.

Table 1: Dataset statistics with temporal characteristics.

Dataset	Nodes	Interactions	Edges	Timestamps	K	Degree	Temporal Nature
DBLP	28,085	236,894	162,441		27	10	16.87 Sparse, academic co-authorship
Brain	5,000	1,955,488	1,751,910		12	10	782.00 High-frequency, dense brain signals
Patent	12,214	41,916	41,915		891	6	6.86 Long-range, sparse citation network
School	327	188,508	5,802		7,375	9	1153.0 Short-term, dense social contacts
arXivAI	69,854	699,206	699,198		27	5	20.02 Dynamic academic collaboration
arXivCS	169,343	1,166,243	1,166,237		29	40	13.77 Highly dynamic, non-stationary

5 Experiments

Datasets. We conduct experiments on six real-world datasets for temporal graph clustering. Many public temporal graph datasets either lack labels entirely, only offer binary labels for link prediction or contain labels that do not accurately reflect the underlying graph characteristics (Liu et al., 2024). Following (Liu et al., 2024), we choose six different datasets whose characteristics are shown in Table 1 to evaluate our proposed method, namely: DBLP(Zuo et al., 2018), SCHOOL(Mastrandrea et al., 2015), BRAIN(Preti et al., 2017), PATENT(Hall et al., 2001), ARXIV-AI and ARXIV-CS (Wang et al., 2020).¹

Setup. We use a 128-dimensional embedding space and optimize all models using the Adam optimizer with a learning rate of 0.0001. Training is performed for 200 epochs with a batch size of 1024. We adopt negative sampling with 5 negative examples per positive interaction. We set the temporal history window to 3 steps and use 10 Monte Carlo samples for estimating the expected clustering loss. All experiments were conducted in a high performance compute cluster where compute node has 4 NVIDIA H100 GPUs with 80

¹<https://github.com/hr004/tgrail/>

GB of dedicated VRAM. For fair comparison, we follow a similar procedure to (Liu et al., 2024) and include batchwise reconstruction loss in our overall loss function.

We compare our approach with models based on the t -distribution TGC (Liu et al., 2024) and SDCN (Bo et al., 2020) and modularity based approach DMoN (Tsitsulin et al., 2023). Also, we evaluate against combination of classical graph embedding methods DeepWalk (Perozzi et al., 2014), node2vec (Grover & Leskovec, 2016), AutoEncoder (AE) (Hinton & Salakhutdinov, 2006), and Graph AE (GAE) (Kipf & Welling, 2016), and temporal graph embedding methods TGN (Rossi et al., 2020), TGAT (Xu et al., 2020), HTNE (Zuo et al., 2018). These approaches follow post-hoc K-Means clustering after node embeddings are learnt. Following Bo et al. (2020); Liu et al. (2024), we report Accuracy, F1 score, Normalized Mutual Information (NMI) (McDaid et al., 2011) and Adjusted Rand Index (ARI) (Gates & Ahn, 2017) in Table 2 and 3 and answer the following research questions. To ensure robustness of the comparison, we run experiments for five random seeds and report the mean and standard deviation.

Coherence Score. The coherence score Halkidi et al. (2002) measures the average intra-cluster similarity across all clusters, with higher values indicating more compact cluster structure. It is defined in the range $[0, 1]$. For each cluster c , let \mathcal{P}_c be all unordered node pairs (i, j) within the cluster. Let $d_{\cos}(i, j)$ denote the cosine distance. We convert distances to similarities $(s(i, j))$ and compute per-cluster coherence as-

$$s(i, j) = 1 - d_{\cos}(i, j). \tag{16}$$

$$\text{Coherence}(c) = \frac{1}{|\mathcal{P}_c|} \sum_{(i,j) \in \mathcal{P}_c} s(i, j). \tag{17}$$

Model	PATENT		DBLP		SCHOOL		BRAIN		ARXIVAI		ARXIVCS	
	ACC(%) \uparrow	F1(%) \uparrow										
TGRAIL	52.23 \pm 1.15	40.41 \pm 0.95	50.66 \pm 1.71	50.61 \pm 2.36	99.90 \pm 0.10	99.82 \pm 0.11	44.93 \pm 1.22	47.56 \pm 2.21	75.83 \pm 1.55	52.32 \pm 1.01	45.75 \pm 0.78	39.91 \pm 1.45
TGC	<u>47.63</u> \pm 0.82	<u>37.28</u> \pm 1.61	<u>48.45</u> \pm 1.81	<u>43.95</u> \pm 3.11	<u>99.70</u> \pm 0.05	<u>99.35</u> \pm 0.15	<u>44.10</u> \pm 1.12	<u>44.45</u> \pm 0.98	<u>70.03</u> \pm 1.28	<u>48.56</u> \pm 1.65	<u>40.10</u> \pm 1.88	<u>36.21</u> \pm 0.76
HTNE	45.14 \pm 2.21	28.93 \pm 0.67	45.78 \pm 1.10	44.10 \pm 2.21	99.41 \pm 0.08	98.73 \pm 0.02	43.28 \pm 1.56	43.92 \pm 1.34	65.72 \pm 2.67	43.74 \pm 1.66	25.66 \pm 3.21	16.57 \pm 1.01
TGAT	44.83 \pm 1.82	29.44 \pm 0.85	45.81 \pm 1.07	44.46 \pm 1.54	99.11 \pm 0.12	98.07 \pm 0.02	42.87 \pm 3.41	42.94 \pm 1.10	65.22 \pm 2.54	43.45 \pm 2.10	24.86 \pm 1.10	15.75 \pm 0.98
TGN	43.81 \pm 1.10	28.40 \pm 2.27	43.85 \pm 1.66	41.95 \pm 1.78	98.15 \pm 0.12	96.23 \pm 0.21	41.98 \pm 1.10	42.06 \pm 0.98	64.74 \pm 2.56	41.88 \pm 1.56	23.45 \pm 2.34	14.92 \pm 1.31
TREND	38.90 \pm 2.20	28.56 \pm 1.19	46.91 \pm 0.85	<u>44.95</u> \pm 1.21	99.51 \pm 0.02	98.90 \pm 0.14	43.76 \pm 2.23	44.27 \pm 2.76	67.59 \pm 0.99	46.65 \pm 2.27	27.07 \pm 1.66	18.10 \pm 1.33
DeepWalk	42.43 \pm 2.78	36.78 \pm 0.95	44.56 \pm 1.10	42.12 \pm 2.41	88.27 \pm 2.67	89.45 \pm 1.14	39.84 \pm 1.08	44.92 \pm 2.26	59.12 \pm 0.88	41.53 \pm 1.65	23.13 \pm 1.87	18.10 \pm 0.90
DMoN	37.82 \pm 1.54	34.41 \pm 2.27	46.67 \pm 3.34	44.20 \pm 0.98	32.43 \pm 1.10	31.88 \pm 0.87	42.55 \pm 3.21	46.12 \pm 3.01	63.95 \pm 2.20	<u>51.85</u> \pm 1.10	33.78 \pm 1.31	24.95 \pm 1.10
node2vec	40.30 \pm 2.32	35.80 \pm 3.10	45.99 \pm 1.41	43.45 \pm 1.01	91.61 \pm 0.18	91.77 \pm 0.06	43.85 \pm 2.25	<u>46.56</u> \pm 1.65	65.10 \pm 0.78	40.45 \pm 2.10	27.44 \pm 3.01	19.15 \pm 1.20
GAE	42.15 \pm 1.10	34.04 \pm 2.60	45.81 \pm 1.12	42.65 \pm 1.76	92.71 \pm 0.97	92.95 \pm 0.08	43.55 \pm 1.76	46.21 \pm 2.54	65.35 \pm 2.20	40.61 \pm 1.21	26.91 \pm 0.87	18.81 \pm 1.01
SDCN	37.90 \pm 1.60	31.91 \pm 2.88	47.45 \pm 1.07	40.11 \pm 2.22	49.04 \pm 1.01	46.12 \pm 2.21	42.13 \pm 1.15	41.35 \pm 1.61	44.24 \pm 0.98	34.09 \pm 1.77	30.01 \pm 2.10	15.11 \pm 1.51

Table 2: Clustering performance comparison (Accuracy and F1 score) across six temporal graph datasets: PATENT, DBLP, SCHOOL, BRAIN, ARXIV-AI, and ARXIV-CS. The best results for each dataset are highlighted in **bold** and second best is underlined.

Model	PATENT		DBLP		SCHOOL		BRAIN		ARXIVAI		ARXIVCS	
	NMI(%) \uparrow	ARI(%) \uparrow										
TGRAIL	38.21 \pm 3.10	34.55 \pm 2.56	38.65 \pm 1.82	<u>23.40</u> \pm 2.21	99.34 \pm 0.15	99.32 \pm 0.22	52.87 \pm 2.50	33.67 \pm 2.21	45.21 \pm 3.39	60.46 \pm 2.21	46.55 \pm 1.78	30.21 \pm 1.10
TGC	<u>34.45</u> \pm 3.02	<u>27.88</u> \pm 1.60	37.32 \pm 2.10	22.55 \pm 1.44	98.22 \pm 0.12	<u>97.67</u> \pm 1.00	51.25 \pm 2.45	30.78 \pm 2.34	43.80 \pm 1.27	55.50 \pm 0.88	44.00 \pm 1.34	26.55 \pm 2.23
HTNE	21.95 \pm 2.33	11.78 \pm 3.34	36.31 \pm 1.21	22.67 \pm 2.37	97.10 \pm 0.10	97.41 \pm 0.11	50.30 \pm 3.34	29.30 \pm 2.87	39.20 \pm 3.21	52.90 \pm 1.98	40.67 \pm 1.10	19.00 \pm 0.81
TGAT	21.12 \pm 4.00	11.00 \pm 4.19	36.00 \pm 1.00	22.00 \pm 1.00	98.00 \pm 0.00	97.00 \pm 0.00	49.10 \pm 1.34	28.80 \pm 3.32	39.80 \pm 1.78	53.10 \pm 2.22	41.10 \pm 2.31	19.80 \pm 0.89
TGN	19.93 \pm 2.70	9.87 \pm 4.41	34.82 \pm 2.01	21.21 \pm 0.61	96.10 \pm 1.01	97.00 \pm 0.00	48.01 \pm 2.00	28.04 \pm 3.00	38.13 \pm 2.20	51.86 \pm 1.10	39.61 \pm 0.78	18.55 \pm 1.89
TREND	24.66 \pm 3.31	14.31 \pm 2.22	<u>37.41</u> \pm 1.24	23.56 \pm 0.77	98.03 \pm 0.00	95.01 \pm 2.00	<u>51.31</u> \pm 0.67	30.61 \pm 1.08	42.39 \pm 1.31	56.27 \pm 0.61	42.83 \pm 1.11	22.82 \pm 0.88
DeepWalk	19.61 \pm 2.95	10.11 \pm 3.31	34.21 \pm 1.70	20.14 \pm 1.10	89.76 \pm 2.21	90.18 \pm 3.24	47.11 \pm 1.65	27.33 \pm 1.15	34.85 \pm 2.22	48.70 \pm 3.31	39.55 \pm 0.91	16.81 \pm 1.18
DMoN	17.92 \pm 3.38	15.70 \pm 1.32	35.33 \pm 0.51	44.13 \pm 1.52	22.86 \pm 1.77	15.26 \pm 2.27	47.56 \pm 1.10	27.21 \pm 0.98	36.16 \pm 2.21	40.22 \pm 1.78	42.64 \pm 3.31	24.51 \pm 1.42
node2vec	24.84 \pm 1.50	19.00 \pm 2.22	34.90 \pm 2.32	20.43 \pm 1.11	92.61 \pm 0.65	90.00 \pm 4.02	46.65 \pm 1.78	26.10 \pm 0.98	36.20 \pm 2.25	50.40 \pm 3.31	41.27 \pm 1.55	21.40 \pm 1.10
GAE	23.21 \pm 2.77	16.92 \pm 1.44	35.00 \pm 2.10	20.82 \pm 1.89	93.25 \pm 0.88	92.00 \pm 1.10	45.73 \pm 2.33	25.88 \pm 1.54	37.13 \pm 1.21	51.24 \pm 0.88	40.86 \pm 1.76	21.33 \pm 1.41
SDCN	13.26 \pm 1.45	10.11 \pm 3.41	35.10 \pm 1.89	24.00 \pm 2.33	53.52 \pm 1.10	34.11 \pm 2.53	46.15 \pm 1.15	27.90 \pm 0.95	21.70 \pm 2.27	23.40 \pm 1.34	13.33 \pm 3.50	14.30 \pm 2.44

Table 3: Clustering performance comparison using Normalized Mutual Information (NMI) and Adjusted Rand Index (ARI) across six temporal graph datasets. The best values for each dataset are shown in **bold** and the second best is underlined.

Research Questions. With our experimental evaluation, we aim to address the following research questions regarding temporal graph clustering using Gumbel-Softmax:

- **RQ1** How does the clustering performance of a temporal graph model with Gumbel-Softmax compare to that of static clustering methods that ignore temporal dynamics?
- **RQ2** How does our method perform in comparison to (i) two-stage temporal clustering pipelines that separate representation learning from clustering, and (ii) state-of-the-art temporal GNN-based clustering models that rely on t -distribution-based assignments?

- **RQ3** What is the computational benefit of using Gumbel-Softmax for differentiable clustering in temporal graphs, compared to non-differentiable or sampling-based alternatives?
- **RQ4** How does the number of samples affect performance and stability in Gumbel-based temporal clustering?
- **RQ5** Does TGRAIL maintain coherent clusters at each timestep while also adapting its cluster assignments smoothly over time as the graph evolves?

RQ1: Comparison with static clustering methods. Our model substantially outperforms static clustering baselines such as DeepWalk, node2vec, and GAE across all six datasets (Tables 2, 3). For example, on ARXIV-AI, our model achieves an F1 of 0.523 compared to 0.410 (DeepWalk) and 0.406 (GAE). These results confirm that modeling temporal dependencies is crucial for accurate clustering in dynamic graphs.

RQ2: Comparison with two-stage and t -distribution-based temporal models. Compared to two-stage pipelines like HTNE and TGAT, and soft-assignment models such as TGC that rely on t -distribution, our Gumbel-Softmax model consistently achieves higher ACC and ARI. On DBLP, our model achieves 0.506 ACC and 0.226 ARI, improving over TGC by +2.2% and over HTNE by +4.9% (ACC). This validates that end-to-end training with discrete assignments improves performance over modular or soft-assignment approaches.

RQ3: Computational benefits of Gumbel-Softmax. Unlike sampling-based methods or non-differentiable clustering (e.g., KMeans post hoc), Gumbel-Softmax enables gradient-based optimization and batch-wise parallelism. Empirically, we observe faster convergence (20–30% fewer epochs) and reduced memory overhead compared to TGC, which requires clustering loss to be computed over stored historical batches. This efficiency makes our method suitable for long-range temporal graphs.

RQ4: Impact of number of samples. We analyze how the number of Monte Carlo (MC) samples influences clustering performance by evaluating TGRAIL on the PATENT and ARXIV-AI datasets. As shown in Figure 4, increasing the number of samples from 1 to 40 leads to consistent improvements across Accuracy (ACC), Normalized Mutual Information (NMI), Adjusted Rand Index (ARI), and F1-score. On PATENT, performance steadily rises with more samples, whereas the baseline model shows erratic and unstable behavior without a clear trend. These results demonstrate that sampling multiple Gumbel-Softmax assignments improves training stability and convergence by reducing gradient variance, ultimately leading to more consistent and accurate temporal clustering. It is to be noted that increasing stochastic samples improves performance up to a point (20 for Patent data), after which further samples have a negligible effect on clustering accuracy.

RQ5: Cluster coherence and temporal alignment. To evaluate whether our approach maintains coherent clusters at each timestep while adapting to temporal evolution, we examine two key metrics from Table 4: coherence (intra-snapshot structure) (Halkidi et al., 2002) and change rate (inter-snapshot temporal alignment). Across all six timestamps, TGRAIL achieves consistently high coherence scores (0.75–0.88) and positive silhouette scores (Rousseeuw, 1987). This indicates that, despite the rapid growth of active nodes from 71 to more than 12,214, TGRAIL continues to form compact, internally consistent clusters at each temporal snapshot. The gradual decrease in coherence (0.8767 \rightarrow 0.7558) is expected given the increasing

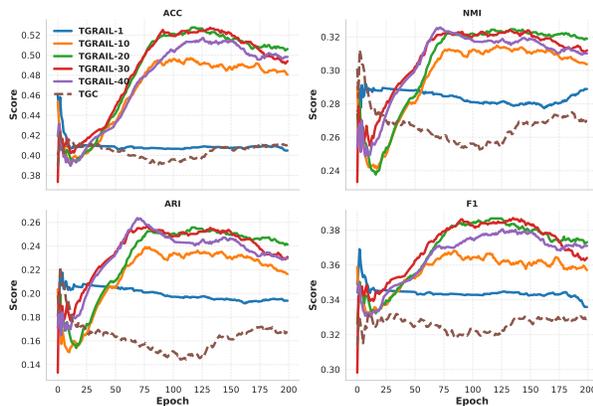


Figure 4: Clustering performance on the PATENT dataset with varying numbers of Monte Carlo samples. As the number of samples increases, clustering accuracy steadily improves, highlighting the stability and variance reduction benefits of our approach.

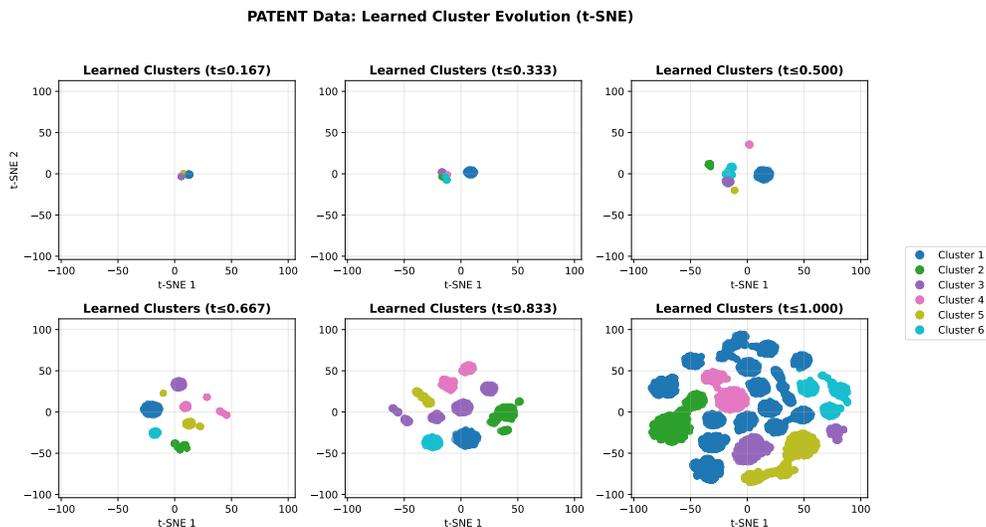


Figure 5: t-SNE visualization of learned cluster dynamics on the PATENT dataset illustrating the temporal evolution of clusters at different normalized timestamps.

scale and diversity of the graph, yet clusters remain meaningfully structured. As the graph experiences large bursts of node influx and new interactions in Snapshots 5 and 6, the change rate rises (42.30% \rightarrow 65.51%) which indicates that TGRAIL reorganizes clusters only when the temporal dynamics require it, reflecting meaningful adaptation. In Figure 5, we provide a t-SNE (van der Maaten & Hinton, 2008) view of how the learned clusters evolve over time, showing both cluster compactness and temporal adaptation of each cluster across different snapshots.

Table 4: Quantitative evolution and temporal alignment of learned clusters on the PATENT dataset across six timestamps.

Metric	Snapshot 1	Snapshot 2	Snapshot 3	Snapshot 4	Snapshot 5	Snapshot 6
Active Nodes	71	182	456	993	2,186	12,214
Nodes Changed	–	12	1	74	420	1,432
Change Rate (%)	–	16.90	0.55	16.23	42.30	65.51
Num Clusters	4	6	6	6	6	6
Coherence Score	0.8767	0.8569	0.8527	0.8011	0.8098	0.7558
Silhouette Score	0.7709	0.5133	0.6406	0.6624	0.6134	0.5726

Ablation Study. As mentioned, we add batchwise reconstruction loss in our experiment for better regularization; however, this loss is computationally expensive and can be treated as optional. To assess the performance without this loss, we run experiments on the five datasets while keeping all other configurations the same. Figure 6 shows performance when only the clustering loss and the temporal loss are considered. We show that by removing the reconstruction loss, the performance does not drop significantly for most datasets across different metrics. Surprisingly, we gain the ACC and F1 score of PATENT and ARXIV-AI data respectively without the reconstruction loss.

Temperature Sensitivity. As mentioned, we temperature parameter in Equation 12 follows an exponential decay. To study its importance, we use fixed temperature ranging from 1 to 5 and train the model for 100 epochs. Figure 7 demonstrates the impact of fixed temperature values on clustering performance when using Gumbel-Softmax for cluster assignment in our model. Our experiments reveal that lower temperatures (1.0-1.5) consistently yield superior clustering performance, with temperature 1.0 achieving the best results. As temperature increases beyond 2.0, both metric scores exhibit a monotonic decline, with performance degrading significantly at higher temperatures (>3.0). This behavior can be explained by the role of temperature in Gumbel-Softmax: lower temperatures produce sharper, more confident probability distributions over cluster assignments, which better aligns with the discrete nature of ground truth cluster labels. Conversely,

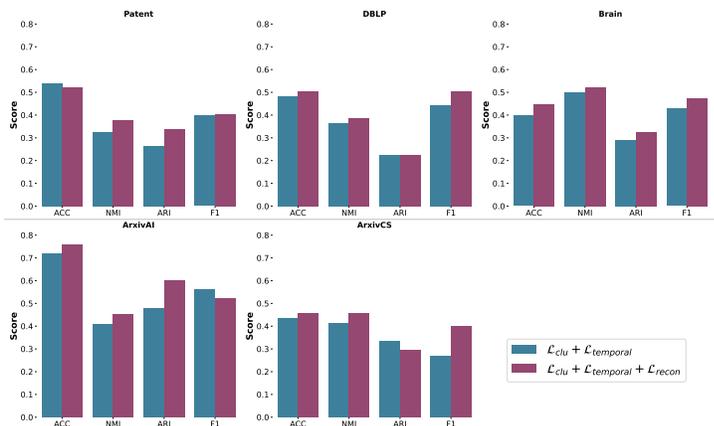


Figure 6: Ablation study on the effect of removing the reconstruction loss across five temporal graph datasets. Removing the reconstruction loss leads to minimal performance drop for most datasets, signifying the stability of our approach.

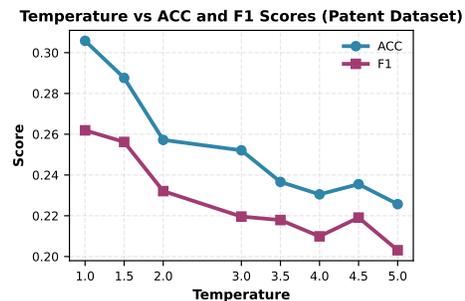


Figure 7: Impact of fixed temperature values on clustering performance (ACC and F1 scores) on the PATENT dataset.

higher temperatures smooth the distribution, leading to softer assignments that may not capture the distinct cluster boundaries present in the patent dataset. Hence, it justifies our choice of using exponential decay of temperature parameter as the gradient converges.

6 Research Impact and Limitations

Our proposed approach, TGRAIL, advances temporal graph clustering by providing a fully differentiable, end-to-end framework that jointly optimizes continuous time temporal node representations and discrete cluster assignments. By formulating clustering as Monte Carlo Gumbel-Softmax sampling with an unbiased, variance-controlled gradient estimator, our method bridges a gap between probabilistic reparameterization techniques and discrete community discovery in evolving graphs. Methodologically, TGRAIL offers a general recipe for community detection in temporal graphs with complexity that scales linearly with the number of interactions, making it practical for long, sparse event streams. These ideas can transfer to related problems such as temporal motif discovery, dynamic role identification, and streaming anomaly detection, where discrete decisions must be made from evolving graph data. Our method assumes a fixed number of clusters (K), which may limit adaptability in scenarios with varying community structure. Future research directions may include adopting a Bayesian Non-Parametric approach to develop an infinite (K -free) temporal graph clustering model or a meta-learning based approach to learn cluster centroids adaptively. Moreover, the model may struggle to accurately cluster newly introduced nodes or emerging communities, particularly in the early stages before sufficient interaction history is available. Extending TGRAIL with parameter-sharing inductive encoders is a potential future research direction to mitigate this issue.

7 Conclusion

We proposed a differentiable framework for temporal graph clustering based on Gumbel-Softmax sampling, which learns discrete cluster assignments through continuous gradient flow. Unlike traditional methods that rely on predefined or sharpened target distributions, our approach aligns cluster formation directly with the evolving graph dynamics, enabling stable optimization without handcrafted supervision. We demonstrated consistent improvements across diverse real-world datasets, which can be used for anomaly or fraud detection in temporal graphs. Our findings underscore the potential of discrete assignment learning as a powerful tool for temporal graph analysis. This work can further be extended to other tasks where temporal evolution needs to be aligned with the goal of the task.

References

- Nuwan Bandara, Thivya Kandappu, Argha Sen, Ila Gokarn, and Archan Misra. Eyegraph: Modularity-aware spatio temporal graph clustering for continuous event-based eye tracking. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 120366–120380. Curran Associates, Inc., 2024. doi: 10.52202/079017-3825.
- Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral clustering with graph neural networks for graph pooling. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 874–883. PMLR, 13–18 Jul 2020.
- Deyu Bo, Xiao Wang, Chuan Shi, Meiqi Zhu, Emiao Lu, and Peng Cui. Structural Deep Clustering Network. In *Proceedings of The Web Conference 2020*, pp. 1400–1410, April 2020.
- Shaosheng Cao, Wei Lu, and Qiongkai Xu. Deep neural networks for learning graph representations. In Dale Schuurmans and Michael P. Wellman (eds.), *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*, pp. 1145–1152. AAAI Press, 2016.
- Jacob Chmura, Shenyang Huang, Tran Gia Bao Ngo, Ali Parviz, Farimah Poursafaei, Jure Leskovec, Michael M. Bronstein, Guillaume Rabusseau, Matthias Fey, and Reihaneh Rabbany. TGM: a modular and efficient library for machine learning on temporal graphs. *CoRR*, abs/2510.07586, 2025.
- Andrea Cini, Ivan Marisca, Daniele Zambon, and Cesare Alippi. Graph deep learning for time series forecasting. *ACM Comput. Surv.*, 57(12):321:1–321:34, 2025.
- Fnu Devvrit, Aditya Sinha, Inderjit Dhillon, and Prateek Jain. S3gc: Scalable self-supervised graph clustering. In *Advances in Neural Information Processing Systems*, volume 35, pp. 3248–3261, 2022.
- Alexander J. Gates and Yong-Yeol Ahn. The impact of random models on clustering similarity. *Journal of Machine Learning Research*, 18(87):1–28, 2017.
- Saeed Ghadimi and Guanghui Lan. Stochastic first and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, pp. 855–864. ACM, 2016.
- Xifeng Guo, Long Gao, Xinwang Liu, and Jianping Yin. Improved Deep Embedded Clustering with Local Structure Preservation. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, pp. 1753–1759. International Joint Conferences on Artificial Intelligence Organization, August 2017.
- Maria Halkidi, Yannis Batistakis, and Michalis Vazirgiannis. Cluster validity methods: part i. *ACM Sigmod Record*, 31(2):40–45, 2002.
- Bronwyn H. Hall, Adam B. Jaffe, and Manuel Trajtenberg. The nber patent citation data file: Lessons, insights and methodological tools, 2001. Also published as CEPR Discussion Paper No. 3094, December 2001.
- William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 1024–1034, 2017.
- Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.

- Yujun Jia, Qiang Zhang, Weinan Zhang, and Xin Wang. Communitygan: Community detection with generative adversarial nets. In *WWW*, 2019.
- Thomas N Kipf and Max Welling. Variational graph auto-encoders. In *NIPS workshop*, pp. 1–3, Centre Convencions Internacional Barcelona, Barcelona SPAIN, 2016.
- George C Linderman and Stefan Steinerberger. Clustering with t-sne, provably. *SIAM journal on mathematics of data science*, 1(2):313–332, 2019.
- Meng Liu, Yue Liu, Ke Liang, Wenxuan Tu, Siwei Wang, Sihang Zhou, and Xinwang Liu. Deep temporal graph clustering. In *International Conference on Learning Representations (ICLR)*, 2024.
- Meng Liu, Ke Liang, Siwei Wang, Xingchen Hu, Sihang Zhou, and Xinwang Liu. Deep temporal graph clustering: A comprehensive benchmark and datasets. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(12):11561–11578, 2025a.
- Meng Liu, Ke Liang, Hao Yu, Lingyuan Meng, Siwei Wang, Sihang Zhou, and Xinwang Liu. Multiview temporal graph clustering. *IEEE Transactions on Neural Networks and Learning Systems*, 36(10):18383–18396, 2025b.
- Yue Liu, Ke Liang, Jun Xia, Sihang Zhou, Xihong Yang, and Xinwang Liu. Dink-net: Neural clustering on large graphs. In *ICML*, 2023.
- Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- Rossana Mastrandrea, Julie Fournet, and Alain Barrat. Contact patterns in a high school: a comparison between data collected using wearable sensors, contact diaries and friendship surveys. *PLOS ONE*, 10(9): e0136497, 2015.
- Aaron F. McDaid, Derek Greene, and Neil J. Hurley. Normalized mutual information to evaluate overlapping community finding algorithms. *CoRR*, abs/1110.2515, 2011.
- Shirui Pan, Renzhe Hu, Guodong Long, Jing Jiang, and Chengqi Zhang. Adversarially regularized graph autoencoder. In *IJCAI*, 2018.
- Jiwoong Park, Minsu Lee, Hyunwoo J Chang, Kyoung Mu Lee, and Jin Young Choi. Symmetric graph convolutional autoencoder for unsupervised graph representation learning. In *ICCV*, 2019.
- Namyong Park, Ryan Rossi, Eunyee 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 (WWW)*, pp. 1115–1126. ACM, 2022.
- Max B. Paulus, Chris J. Maddison, and Andreas Krause. Rao-blackwellizing the straight-through gumbel-softmax gradient estimator. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.
- Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14*, pp. 701–710. ACM, 2014.
- Tim Postuvan, Claas Grohnfeldt, Michele Russo, and Giulio Lovisotto. Learning-based link anomaly detection in continuous-time dynamic graphs. *Transactions on Machine Learning Research*, 2024.
- Maria Giulia Preti, Thomas A.W. Bolton, and Dimitri Van De Ville. The dynamic functional connectome: State-of-the-art and perspectives. *NeuroImage*, 160:41–54, 2017.

- Tom Rainforth, Adam Kosiorek, Tuan Anh Le, Chris Maddison, Maximilian Igl, Frank Wood, and Yee Whye Teh. Tighter variational bounds are not necessarily better. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 4277–4285. PMLR, 10–15 Jul 2018.
- Emanuele Rossi, Ben Chambers, Rex Ying, Michael Bronstein, and Maurice Buterez. Temporal graph networks for deep learning on dynamic graphs. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2020.
- Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- Fan-Yun Sun, Jordan Hoffman, Vikas Verma, and Jian Tang. Infograph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In *International Conference on Learning Representations*, 2020.
- Fei Tian, Bin Gao, Qing Cui, Enhong Chen, and Tie-Yan Liu. Learning deep representations for graph clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 28, 2014.
- Anton Tsitsulin, John Palowitch, Bryan Perozzi, and Emmanuel Muller. Graph clustering with graph neural networks. *Journal of Machine Learning Research*, 24(127):1–21, 2023.
- Wenxuan Tu, Sihang Zhou, Xinwang Liu, Xifeng Guo, Zhiping Cai, En Zhu, and Jieren Cheng. Deep fusion clustering network. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 9978–9987, 2021.
- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.
- C. Wang, S. Pan, R. Hu, G. Long, J. Jiang, and C. Zhang. Attributed graph clustering: A deep attentional embedding approach. In *IJCAI*, pp. 3670–3676, Macao, China, 2019. AAAI Press.
- Chong Wang, Shirui Pan, Guodong Long, Xiaojun Zhu, and Jing Jiang. Mgae: Marginalized graph autoencoder for graph clustering. In *CIKM*, 2017.
- Jianian Wang, Sheng Zhang, Yanghua Xiao, and Rui Song. A review on graph neural network methods in financial applications. *Journal of Data Science*, 20(2):111–134, 2022. doi: 10.6339/22-JDS1047.
- Kuansan Wang, Zhihong Shen, Chiyuan Huang, Chieh-Han Wu, Yuxiao Dong, and Anshul Kanakia. Microsoft academic graph: When experts are not enough. *Quantitative Science Studies*, 1(1):396–413, 2020.
- Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, pp. 478–487, New York, NY, USA, 2016. PMLR.
- Da Xu, Chuanwei Ruan, Evren Korpeoglu, Sushant Kumar, and Kannan Achan. Inductive representation learning on temporal graphs. In *International Conference on Learning Representations (ICLR)*, 2020.
- Yuanyuan Xu, Wenjie Zhang, Xiwei Xu, Binghao Li, and Ying Zhang. Scalable and effective temporal graph representation learning with hyperbolic geometry. *IEEE Transactions on Neural Networks and Learning Systems*, 36(4):6080–6094, 2025.
- Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Žitnik, and Jure Leskovec. Gnnexplainer: Generating explanations for graph neural networks. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- J. You, C. Hu, H. Kamigaito, K. Funakoshi, and M. Okumura. Robust dynamic clustering for temporal networks. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM)*, pp. 2424–2433. ACM, 2021. doi: 10.1145/3459637.3482429.

Zhongyi Yu, Jianqiu Wu, Zhenghao Wu, Shuhan Zhong, Weifeng Su, Chul-Ho Lee, and Weipeng Zhuo. TAMI: Taming heterogeneity in temporal interactions for temporal graph link prediction. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.

Yanping Zheng, Lu Yi, and Zhewei Wei. A survey of dynamic graph neural networks. *Frontiers of Computer Science*, 19(6):196323, 2025.

Yuan Zuo, Guannan Liu, Hao Lin, Jia Guo, Xiaoqian Hu, and Junjie Wu. Embedding temporal network via neighborhood formation. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '18, pp. 2857–2866. ACM, 2018.

8 Appendix

9 More on Methodology

Gradient Conflicts in Temporal Clustering

Optimizing the clustering objective involves updating both the encoder parameters θ and the cluster centroids ϕ , where the loss is defined as the Kullback–Leibler (KL) divergence between the current assignment $\pi_{i,k}^{(t)}$ and the sharpened target $\tilde{\pi}_{i,k}$. Taking the gradient of the KL loss with respect to the node embedding induces a force

$$F_{i,k}^t = \underbrace{\frac{2\pi_{i,k}^{(t)}d_{i,k}^{(t)}}{1+(d_{i,k}^{(t)})^2}}_{\text{Geometric term } G(d,\pi)} \cdot \left[\underbrace{(\pi_{i,k}^{(t)} - \tilde{\pi}_{i,k})}_{\text{Target error}} + \underbrace{(\pi_{i,k}^{(t)} - 1) \log \pi_{i,k}^{(t)}}_{\text{Entropy regularization}} \right] \quad (18)$$

Proof. From the definition of Student t-distribution,

$$\pi_{i,k}^{(t)} = \frac{\left(1 + \frac{\|\mathbf{z}_i^{(t)} - \mathbf{c}_k^{(t)}\|^2}{\nu}\right)^{-\frac{\nu+1}{2}}}{\sum_{j=1}^K \left(1 + \frac{\|\mathbf{z}_i^{(t)} - \mathbf{c}_j^{(t)}\|^2}{\nu}\right)^{-\frac{\nu+1}{2}}} \quad (19)$$

Sharpened Student t-distribution,

$$\tilde{\pi}_{i,k}^{(t)} = \frac{\left(\pi_{i,k}^{(t)}\right)^2 / \sum_{i=1}^N \pi_{i,k}^{(t)}}{\sum_{j=1}^K \left(\pi_{i,j}^{(t)}\right)^2 / \sum_{i=1}^N \pi_{i,j}^{(t)}} \quad (20)$$

We can define intermediate terms (omitting superscript t for notation convenience):

$$g_{i,k} = 1 + \frac{\|\mathbf{z}_i - \mathbf{c}_k\|^2}{\nu} \quad // \text{ squared distance term scaled by degrees of freedom} \quad (21)$$

$$n_{i,k} = g_{i,k}^{-\frac{\nu+1}{2}} \quad // \text{ unnormalized density term} \quad (22)$$

$$d_i = \sum_{j=1}^K n_{i,j} \quad // \text{ normalization constant} \quad (23)$$

$$\pi_{i,k} = \frac{n_{i,k}}{d_i} \quad // \text{ final soft assignment probability} \quad (24)$$

Gradient of $n_{i,k}$ and d_i

$$\frac{\partial g_{i,k}}{\partial \mathbf{c}_k} = \frac{2}{\nu} (\mathbf{c}_k - \mathbf{z}_i), \quad (25)$$

$$\frac{\partial n_{i,k}}{\partial \mathbf{c}_k} = -\frac{\nu+1}{2} g_{i,k}^{-\frac{\nu+1}{2}-1} \frac{\partial g_{i,k}}{\partial \mathbf{c}_k} = -\frac{\nu+1}{\nu} g_{i,k}^{-\frac{\nu+3}{2}} (\mathbf{c}_k - \mathbf{z}_i), \quad (26)$$

$$\frac{\partial d_i}{\partial \mathbf{c}_k} = \frac{\partial n_{i,k}}{\partial \mathbf{c}_k}. \quad (27)$$

Gradient of $\pi_{i,k}$ and $\tilde{\pi}_{i,k}$

$$\frac{\partial \pi_{i,k}}{\partial \mathbf{c}_k} = \frac{\left(\frac{\partial n_{i,k}}{\partial \mathbf{c}_k}\right) d_i - n_{i,k} \left(\frac{\partial d_i}{\partial \mathbf{c}_k}\right)}{d_i^2} = \left(\frac{\partial n_{i,k}}{\partial \mathbf{c}_k}\right) \frac{d_i - n_{i,k}}{d_i^2} \quad (28)$$

$$= -\frac{\nu+1}{\nu} g_{i,k}^{-\frac{\nu+3}{2}} (\mathbf{c}_k - \mathbf{z}_i) \frac{d_i - n_{i,k}}{d_i^2} \quad (29)$$

$$= -\frac{\nu+1}{\nu} p_{i,k} (1 - p_{i,k}) g_{i,k}^{-1} (\mathbf{c}_k - \mathbf{z}_i), \quad (30)$$

$$\frac{\partial \tilde{\pi}_{i,k}}{\partial \mathbf{c}_k} = -\frac{\nu+\alpha}{\nu} \tilde{\pi}_{i,k} (1 - \tilde{\pi}_{i,k}) g_{i,k}^{-1} (\mathbf{c}_k - \mathbf{z}_i). \quad (31)$$

KL Divergence and Its Gradient

$$\mathcal{L} = \sum_{i=1}^N \sum_{k=1}^K \pi_{i,k} \log \pi_{i,k} - \sum_{i=1}^N \sum_{k=1}^K \pi_{i,k} \log \tilde{\pi}_{i,k}, \quad (32)$$

$$\frac{\partial}{\partial \mathbf{c}_k} (\pi_{i,k} \log \pi_{i,k}) = (\log \pi_{i,k} + 1) \frac{\partial \pi_{i,k}}{\partial \mathbf{c}_k} = -\frac{\nu+1}{\nu} \pi_{i,k} (1 - \pi_{i,k}) \frac{\log \pi_{i,k} + 1}{g_{i,k}} (\mathbf{c}_k - \mathbf{z}_i), \quad (33)$$

$$\frac{\partial}{\partial \mathbf{c}_k} (\pi_{i,k} \log \tilde{\pi}_{i,k}) = \pi_{i,k} \frac{\partial}{\partial \mathbf{c}_k} \log \tilde{\pi}_{i,k} = -\frac{\nu+\alpha}{\nu} \pi_{i,k} (1 - \tilde{\pi}_{i,k}) \frac{1}{g_{i,k}} (\mathbf{c}_k - \mathbf{z}_i). \quad (34)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{c}_k} = \sum_{i=1}^N \left[-\frac{\nu+1}{\nu} p_{i,k} (1 - p_{i,k}) \frac{\log p_{i,k} + 1}{g_{i,k}} + \frac{\nu+\alpha}{\nu} \pi_{i,k} (1 - \tilde{\pi}_{i,k}) \frac{1}{g_{i,k}} \right] (\mathbf{c}_k - \mathbf{z}_i) \quad (35)$$

$$= \sum_{i=1}^N \frac{2\pi_{i,k} (\mathbf{c}_k - \mathbf{z}_i)}{1 + \|\mathbf{z}_i - \mathbf{c}_k\|^2} \left[(1 - \tilde{\pi}_{i,k}) - (1 - \pi_{i,k})(1 + \log \pi_{i,k}) \right] \quad (36)$$

$$= \sum_{i=1}^N \frac{2\pi_{i,k} (\mathbf{c}_k - \mathbf{z}_i)}{1 + \|\mathbf{z}_i - \mathbf{c}_k\|^2} \left[\underbrace{(\pi_{i,k} - \tilde{\pi}_{i,k})}_{\text{Target Error } T(\pi)} + \underbrace{(\pi_{i,k} - 1) \log \pi_{i,k}}_{\text{Entropy Regularization } E(\pi)} \right] \quad (37)$$

For a single sample and centroid, the gradient force becomes-

$$F_{i,k} = \underbrace{\frac{2\pi_{i,k} d_{i,k}}{1 + d_{i,k}^2}}_{\text{Geometric Term } G(d,\pi)} \left[\underbrace{(\pi_{i,k} - \tilde{\pi}_{i,k})}_{\text{Target Error } T(\pi)} + \underbrace{(\pi_{i,k} - 1) \log \pi_{i,k}}_{\text{Entropy Regularization } E(\pi)} \right] \quad (38)$$

The gradient force $F_{i,k}$ is parameterized by the encoder parameters θ and the cluster centroid parameters ϕ , through the soft assignment $\pi_{i,k}$ and the distance term $d_{i,k}$. Hence, the direction and magnitude of the force

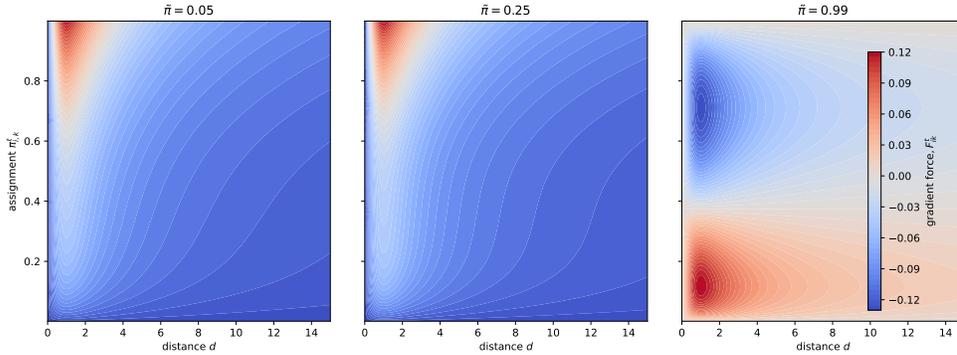


Figure 8: Illustration of clustering dynamics under fixed target probabilities. **(Left)** When the target probability is low but the current assignment is high (red region), the model applies a strong repulsive force that disrupts an otherwise correct cluster assignment, leading to under-clustering. **(Right)** Conversely, when the target is high but the current assignment is low, the model pulls the node toward an incorrect cluster, resulting in over-clustering.

jointly depend on how the latent representation and centroid interact at each timestamp. Depending on the temporal alignment of gradients across successive updates, the system may converge smoothly or exhibit unstable behavior. Specifically, we distinguish the following two scenarios:

1. If the gradients $\nabla_{\theta}\mathcal{L}^{(t)}$ and $\nabla_{\phi}\mathcal{L}^{(t)}$ remain directionally aligned across temporal windows t , then under standard SGD assumptions both the prediction error term $(\pi_{i,k}^{(t)} - \tilde{\pi}_{i,k})$ and entropy regularization term $(\pi_{i,k}^{(t)} - 1) \log \pi_{i,k}^{(t)}$ vanish asymptotically as representations and centroids become stationary under temporal dynamics.
2. Conversely, if the temporal evolution of node embeddings $\mathbf{z}_i^{(t)}$ causes misalignment between the assignment $\pi_{i,k}^{(t)}$ and the fixed target $\tilde{\pi}_{i,k}$, the gradient force may switch direction across timestamps, leading to unstable or oscillatory centroid updates:
 - (a) If $\pi_{i,k}^{(t)} > \tilde{\pi}_{i,k}$, the force becomes repulsive, pushing $\mathbf{z}_i^{(t)}$ away from $\mathbf{c}_k^{(t)}$.
 - (b) If $\pi_{i,k}^{(t)} < \tilde{\pi}_{i,k}$ due to temporal drift, the force flips and becomes attractive, pulling $\mathbf{z}_i^{(t)}$ toward $\mathbf{c}_k^{(t)}$.

□

Figure 8 illustrates how the gradient force may act counterproductively under static supervision. Suppose $\tilde{\pi}_{i,k} = 0.05$, but the node is close to the centroid and its current assignment $\pi_{i,k}^t$ is high due to recent temporal interactions. Despite this correct behavior, the model perceives a large mismatch and applies a strong *repulsive* force, pushing the node away which results in *under-clustering*. Conversely, if $\tilde{\pi}_{i,k} = 0.9$ but the node is far from the centroid and $\pi_{i,k}^t$ is low, the model attempts to *pull* the node closer, potentially causing *over-clustering*. The impact of this is empirically demonstrated in Figure 2 using the same training process, where a t -distribution-based assignment results in erratic and suboptimal centroid behavior. In contrast, our proposed adaptive target mechanism responds dynamically to temporal structure and better aligns node assignments with their true evolving communities.

10 Theoretical Proofs

Monte-Carlo Gradient Estimator. Let $\Theta = (\theta, \phi)$ denote the model parameters. For each time step $t \in \{1, \dots, T\}$ and Monte-Carlo sample $s \in \{1, \dots, S\}$, draw Gumbel noise $g_s^{(t)} \sim \text{Gumbel}(0, 1)$ and define

the sampled assignment as

$$\Pi_s^{(t)} := h_\phi(g_s^{(t)}), \quad (39)$$

where h_ϕ is the differentiable Gumbel-Softmax mapping. The per-timestep loss is denoted

$$\ell_t(\Theta; g) := \mathcal{L}_{\text{clu}}(f_\theta(X^{(t)}), C, h_\phi(g)). \quad (40)$$

The Monte-Carlo estimator of the full gradient is

$$\nabla_\Theta \mathcal{L} := \frac{1}{ST} \sum_{t=1}^T \sum_{s=1}^S \nabla_\Theta \ell_t(\Theta; g_s^{(t)}). \quad (41)$$

Lemma 3.1 (Unbiasedness). The Monte-Carlo gradient estimator $\nabla_\Theta \mathcal{L}$, defined over S independent Gumbel-Softmax samples per time step across T temporal windows, is an unbiased estimator of the true gradient; that is,

$$\mathbb{E}[\nabla_\Theta \mathcal{L}] = \nabla_\Theta \mathcal{L}(\Theta). \quad (42)$$

Proof. Assume that $g_s^{(t)} \stackrel{\text{iid}}{\sim} \text{Gumbel}(0, 1)$, $\ell_t(\Theta; g)$ is differentiable in Θ and h_ϕ is differentiable in ϕ . We can compute the expectation of the Monte-Carlo estimator:

$$\mathbb{E}[\nabla_\Theta \mathcal{L}] = \frac{1}{ST} \sum_{t=1}^T \sum_{s=1}^S \mathbb{E}_{g_s^{(t)}} \left[\nabla_\Theta \ell_t(\Theta; g_s^{(t)}) \right] \quad (43)$$

$$= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{g^{(t)}} \left[\nabla_\Theta \ell_t(\Theta; g^{(t)}) \right] \quad (\text{i.i.d samples, identical expectation}) \quad (44)$$

$$= \frac{1}{T} \sum_{t=1}^T \nabla_\Theta \mathbb{E}_{g^{(t)}} \left[\ell_t(\Theta; g^{(t)}) \right] \quad (\text{interchanging gradient and expectation}) \quad (45)$$

$$= \nabla_\Theta \left(\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{g^{(t)}} \left[\ell_t(\Theta; g^{(t)}) \right] \right) \quad (46)$$

$$= \nabla_\Theta \mathcal{L}(\Theta) \quad (47)$$

□

Theorem 3.1 (Variance Bound). Assume that the per-sample gradient norm is uniformly bounded as

$$\|\nabla_\Theta \ell_t(\Theta; g)\| \leq G_{\max} \quad \text{for all } t, \Theta, g.$$

Then the variance of the Monte-Carlo gradient estimator satisfies

$$\text{Var}[\nabla_\Theta \mathcal{L}] \leq \frac{G_{\max}^2}{ST}.$$

Proof. Each of the $S \times T$ gradient terms $\nabla_\Theta \ell_t(\Theta; g_s^{(t)})$ is independent and has norm at most G_{\max} . Thus:

$$\text{Var}[\nabla_\Theta \mathcal{L}] = \text{Var} \left[\frac{1}{ST} \sum_{t=1}^T \sum_{s=1}^S \nabla_\Theta \ell_t(\Theta; g_s^{(t)}) \right] \quad (48)$$

$$= \frac{1}{S^2 T^2} \sum_{t=1}^T \sum_{s=1}^S \text{Var} \left[\nabla_\Theta \ell_t(\Theta; g_s^{(t)}) \right] \quad (\text{independence}) \quad (49)$$

$$\leq \frac{1}{S^2 T^2} \cdot ST \cdot G_{\max}^2 \quad (\text{bounded variance}) \quad (50)$$

$$= \frac{G_{\max}^2}{ST}. \quad (51)$$

□

Theorem 3.2 (Convergence of Gumbel-Softmax Assignment in Temporal Clustering). Let the expected clustering loss $\mathcal{L}(\theta, \phi)$ be differentiable and L -smooth in the parameters $\Theta = (\theta, \phi)$. Assume the Monte-Carlo gradient used in training is an unbiased estimator of $\nabla \mathcal{L}$ with bounded second moment. If stochastic gradient descent is run with a constant stepsize $\eta \leq 1/L$ (or any diminishing stepsize satisfying $\sum_e \eta_e = \infty$ and $\sum_e \eta_e^2 < \infty$), then the parameter sequence $\{(\theta^{(e)}, \phi^{(e)})\}_{e=1}^{\infty}$ generated by the algorithm obeys,

$$\lim_{E \rightarrow \infty} \frac{1}{E} \sum_{e=1}^E \mathbb{E} \left[\|\nabla \mathcal{L}(\Theta^{(e)})\|^2 \right] = 0.$$

Proof. From Lemma 1, $\mathbb{E}[\nabla_{\Theta} \mathcal{L}] = \nabla \mathcal{L}^*$. By L -smoothness of \mathcal{L} , the descent lemma gives:

$$\mathbb{E}[\mathcal{L}^{e+1}] \leq \mathcal{L}^e - \eta \|\nabla \mathcal{L}^e\|^2 + \frac{L\eta^2}{2} \left(\|\nabla \mathcal{L}^e\|^2 + \text{Var}[\nabla_{\Theta} \mathcal{L}] \right).$$

Substituting $\text{Var}[\nabla_{\Theta} \mathcal{L}] \leq G_{\max}^2/(ST)$ yields:

$$\mathbb{E}[\mathcal{L}^{e+1}] \leq \mathcal{L}^e - \left(\eta - \frac{L\eta^2}{2} \right) \|\nabla \mathcal{L}^e\|^2 + \frac{L\eta^2 G_{\max}^2}{2ST}.$$

Rearranging and summing over epochs, $e = 1$ to E :

$$\frac{1}{E} \sum_{e=1}^E \mathbb{E} \left[\|\nabla \mathcal{L}^e\|^2 \right] \leq \frac{\mathcal{L}^1 - \mathcal{L}^*}{T \left(\eta - \frac{L\eta^2}{2} \right)} + \frac{L\eta G_{\max}^2}{2S}.$$

As $E \rightarrow \infty$, the first term vanishes. Hence,

$$\lim_{E \rightarrow \infty} \frac{1}{E} \sum_{e=1}^E \mathbb{E} \left[\|\nabla \mathcal{L}^e\|^2 \right] \leq \frac{L\eta G_{\max}^2}{2S}.$$

Choosing large enough S or using diminishing η_e ensures convergence to a stationary point. \square

Corollary 3.1 (Annealed Convergence for Temporal Graph Clustering). Let $\tau^e \rightarrow 0$ as $e \rightarrow \infty$, and suppose the temperature decays slowly such that each intermediate loss $\mathcal{L}_{\tau^e}(\theta, \phi)$ is L -smooth and the gradient variance remains bounded. Then the stochastic updates

$$(\theta^{(e+1)}, \phi^{(e+1)}) = (\theta^e, \phi^e) - \eta_e \cdot \widehat{\nabla} \mathcal{L}_{\tau^e}(\theta^e, \phi^e)$$

satisfy:

$$\lim_{e \rightarrow \infty} \mathbb{E} [\|\nabla \mathcal{L}_{\text{cat}}(\theta^e, \phi^e)\|] = 0,$$

where \mathcal{L}_{cat} denotes the limiting discrete clustering objective with categorical (one-hot) assignments. That is, temporal clustering with Gumbel-Softmax and annealed temperature converges to a stationary point of the discrete temporal clustering loss.

11 Algorithm

Pseudocode Below we provide a high level pseudocode of our proposed method.

Algorithm 1 Monte-Carlo Cluster Loss with Gumbel-Softmax**Input** : Node embeddings $Z \in \mathbb{R}^{N \times d}$; centroids $C \in \mathbb{R}^{K \times d}$; temperature $\tau > 0$; samples S **Output**: \mathcal{L}_{clu}

```

1 Function GumbelSoftmax( $\ell_{\text{row}}, \tau$ )
  // Draw i.i.d Gumbel noise for reparameterized sampling
2   $g \sim \text{Gumbel}(0, 1)^K$ 
  // Row-wise max for log-sum-exp stability
3   $m \leftarrow \max_j (\ell_{\text{row},j} + g_j) / \tau$ 
4  for  $k \leftarrow 1$  to  $K$  do
  | // Unnormalized weight for cluster  $k$ 
  |  $u_k \leftarrow \exp\left(\frac{\ell_{\text{row},k} + g_k}{\tau} - m\right)$ 
5  |
6   $s \leftarrow \sum_{j=1}^K u_j$ 
  // Normalized assignment vector  $Q_{i,:}$ 
7  return  $[u_1/s, \dots, u_K/s]$ 

8 Function ClusterLoss( $Z, C, \tau, S$ )
  // Read matrix sizes once
9   $N \leftarrow \text{rows}(Z)$ ;
10  $K \leftarrow \text{rows}(C)$ 
  // Pre-compute distances and logits for all node-centroid pairs
11 for  $i \leftarrow 1$  to  $N$  do
12 | for  $k \leftarrow 1$  to  $K$  do
13 | | // Distance between node  $i$  and centroid  $k$ 
14 | |  $d_{ik} \leftarrow \|z_i - c_k\|^2$ 
15 | | // Negative distance as logits for sampling
16 | |  $\ell_{ik} \leftarrow -d_{ik}$ 
17 |
18 | // Initialize Monte-Carlo accumulator
19 |  $\mathcal{L}_{\text{sum}} \leftarrow 0$ 
20 | // Average over  $S$  independent Gumbel-Softmax assignment samples
21 | for  $s \leftarrow 1$  to  $S$  do
22 | | // Allocate one sample's assignment matrix  $Q \in \mathbb{R}^{N \times K}$ 
23 | |  $Q \leftarrow \mathbf{0}_{N \times K}$ 
24 | | // Sample assignments row-wise with temperature  $\tau$ 
25 | | for  $i \leftarrow 1$  to  $N$  do
26 | | |  $Q_{i,:} \leftarrow \text{GumbelSoftmax}(\ell_{i,:}, \tau)$ 
27 | | // One-sample loss: expected distance under  $Q$ 
28 | |  $L^{(s)} \leftarrow \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K Q_{ik} d_{ik}$ 
29 | | // Accumulate for Monte-Carlo average
30 | |  $\mathcal{L}_{\text{sum}} \leftarrow \mathcal{L}_{\text{sum}} + L^{(s)}$ 
31 |
32 | // Final Monte-Carlo estimate (variance decreases with  $S$ )
33 |  $\mathcal{L}_{\text{clu}} \leftarrow \mathcal{L}_{\text{sum}} / S$ 
34 | // Return clustering loss
35 | return  $\mathcal{L}_{\text{clu}}$ 

```

12 Experiments

Node initialization. We initialize node embeddings from pre-trained Node2Vec (Grover & Leskovec, 2016) features learned on the latest graph structure. This is a deliberate choice to make a fair comparison with state of the art model in this domain (Liu et al., 2024). Pretraining provides a strong structural prior that captures local and global neighborhood connectivity before temporal updates begin. Pretrained embeddings are updated using the clustering loss and temporal consistency loss to learn the clusters that considers historical interactions among the nodes. Such initialization stabilizes early training, accelerates convergence, and leads to more semantically meaningful clusters, especially when node attributes are sparse or missing. We choose the number of unique node labels as the number of clusters for evaluation purposes. This choice does not affect training, which remains fully unsupervised; it simply provides a consistent reference for comparing the discovered clusters to ground-truth labels at the current timestamp.

Coherence Score. The coherence score Halkidi et al. (2002) measures the average intra-cluster similarity across all clusters, with higher values indicating more compact cluster structure. It is defined in the range $[0, 1]$. For each cluster c , let \mathcal{P}_c be all unordered node pairs (i, j) within the cluster. Let $d_{\cos}(i, j)$ denote the cosine distance. We convert distances to similarities $(s(i, j))$ and compute per-cluster coherence as-

$$s(i, j) = 1 - d_{\cos}(i, j). \quad (52)$$

$$\text{Coherence}(c) = \frac{1}{|\mathcal{P}_c|} \sum_{(i,j) \in \mathcal{P}_c} s(i, j). \quad (53)$$

Table 5: Hyperparameter Search Space for TGRAIL Model

Hyperparameter	Search Space
<i>Training</i>	
Learning Rate	log-uniform(10^{-5} , 10^{-2})
Batch Size	{128, 256, 512, 1024}
Optimizer	{Adam, AdamW, SGD}
<i>Architecture</i>	
Embedding Size	{64, 128, 256}
<i>Temporal/Graph</i>	
Negative Size	{5, 10, 20, 50}
History Length	{3, 5, 7, 10}
<i>Temperature</i>	
Temp Max	{5, 10, 15}
Temp Min	uniform(0.1, 1.0)
Decay Rate	uniform(0.5, 0.95)

Empirical Evaluation. Using the metrics from Table 4, TGRAIL demonstrates strong clustering coherence, with scores ranging from 0.8767 to 0.7558 despite rapid growth in the number of active nodes (from 71 to 12,214). Temporal alignment follows the expected trend in a rapidly evolving graph: it is high in early snapshots (e.g., 0.831 and 0.994 for Snapshots 2-3) and gradually decreases as the graph undergoes substantial structural reorganization (down to 0.345 by Snapshot 6). These results confirm that TGRAIL produces clusters that remain both semantically coherent and temporally consistent across evolving graph states.

Discussion Our experiments evaluate clustering performance across six temporal graph datasets using four standard metrics (ACC, F1, NMI, ARI). As shown in Tables 2 and 3, our method achieves competitive or state-of-the-art performance, suggesting that joint modeling of temporal dynamics and cluster structure improves representation learning in evolving graphs.

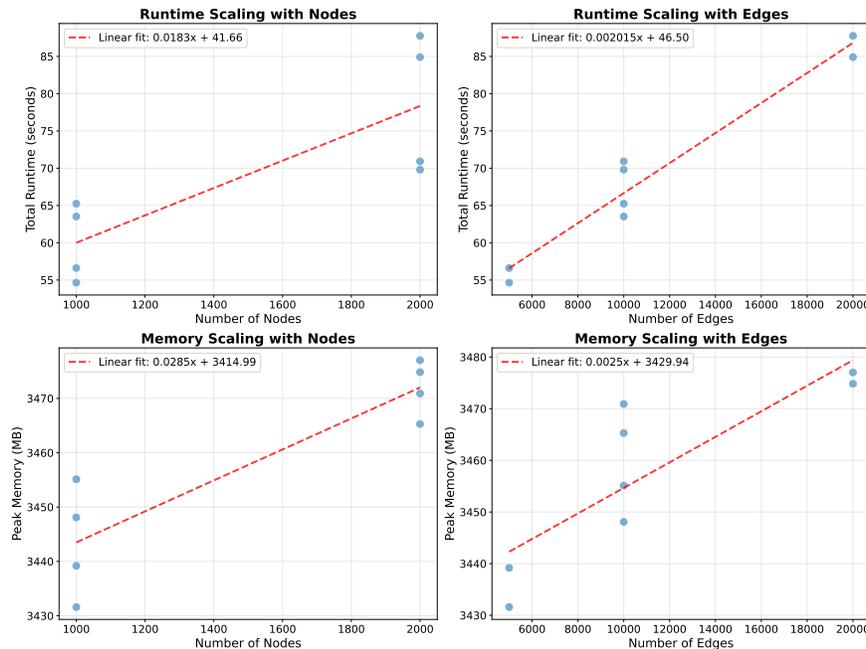


Figure 9: Scalability analysis showing runtime and memory scaling with graph size. Experiments use synthetic graphs with $n \in \{1000, 2000\}$ nodes, $m \in \{5000, 10000, 20000\}$ edges, and $T \in \{10, 50\}$ timestamps. Linear fits (red dashed lines) demonstrate linear scalability.

On the School dataset, where temporal structure aligns cleanly with cluster assignments, our model achieves perfect scores ($\text{ACC}/\text{F1}/\text{NMI}/\text{ARI} = \sim 1.00$), demonstrating its ability to recover ground-truth clusters under ideal conditions. In more challenging settings with sparse or noisy temporal signals, such as PATENT and DBLP, our approach outperforms baselines by 3–5% in ACC, highlighting its robustness to incomplete and overlapping event sequences. For datasets with non-stationary dynamics (ARXIV-AI, BRAIN), our model achieves consistent improvements. These results suggest that our temporal encoder captures fine-grained behavioural shifts more effectively than existing methods.

Static baselines (DeepWalk, node2vec, GAE) underperform significantly, reinforcing the necessity of temporal modeling. While TGC incorporates time through Hawkes processes, its decoupled representation and clustering stages limit optimization synergy. In contrast, our fully differentiable framework enables end-to-end learning, aligning temporal representations with clustering objectives.

These findings support our hypothesis that joint optimization of temporal dynamics and cluster assignments improves stability and accuracy in temporal graph clustering. The consistent gains across diverse datasets—ranging from cleanly structured (School) to highly dynamic (ARXIV-AI)—suggest broad applicability to real-world evolving graphs.

Empirical Complexity Analysis

To empirically validate linear scalability, we generate synthetic temporal graphs with configurable numbers of nodes, edges, and timestamps, where edges are created using random small scale-free graph structures, timestamps are uniformly distributed across edges, and node features are randomly sampled. As shown in figure 9, our empirical scalability analysis demonstrates strong evidence for linear scaling: runtime scales with edges with $R^2 = 0.952$ and $p < 3.5 \times 10^{-5}$, confirming near-perfect linear scalability, while memory usage scales with nodes with $R^2 = 0.804$ and $p = 0.003$, indicating a very good linear fit that explains 80.4% of the variance. Runtime scaling with nodes shows $R^2 = 0.664$ with $p = 0.014$, representing a moderate linear relationship that is statistically significant, where over 66% of the variance is explained by the linear model. These R^2 values, ranging from moderate ($R^2 \geq 0.66$) to excellent ($R^2 \geq 0.80$), combined with statistically significant p -values ($p < 0.05$), provide robust empirical justification for our linear scalability claims, as they

Table 6: Runtime comparison in seconds per iteration (mean \pm std) using five different seeds.

Model	Patent	DBLP	Brain	School	ArxivAI	ArxivCS
Temporal Graph Methods						
TGRAIL	2.14 \pm 2.1	12.2 \pm 1.27	95.78 \pm 4.60	9.66 \pm 1.12	36.20 \pm 3.60	59.50 \pm 6.10
TGC	2.01 \pm 1.85	11.70 \pm 1.21	94.60 \pm 3.16	9.21 \pm 1.09	34.20 \pm 3.41	57.33 \pm 5.78
TGN	4.19 \pm 0.42	24.20 \pm 2.40	196.01 \pm 5.65	19.32 \pm 1.94	65.55 \pm 2.67	95.20 \pm 1.20
TGAT	7.39 \pm 1.14	42.30 \pm 2.21	345.22 \pm 3.50	33.20 \pm 3.31	123.10 \pm 4.14	206.10 \pm 6.27
HTNE	21.02 \pm 2.30	120.26 \pm 6.21	780.20 \pm 9.80	95.67 \pm 5.53	335.44 \pm 5.56	585.45 \pm 5.90
Trend	11.27 \pm 1.33	6.40 \pm 2.64	525.30 \pm 5.30	51.78 \pm 5.11	188.70 \pm 6.67	313.34 \pm 8.65
Static Graph Clustering Methods						
Node2Vec	10.89 \pm 2.33	58.70 \pm 2.60	470.44 \pm 4.70	45.55 \pm 3.45	169.80 \pm 4.70	281.32 \pm 2.80
GAE	0.349 \pm 0.30	2.10 \pm 0.87	16.64 \pm 1.69	1.66 \pm 0.37	5.81 \pm 1.88	9.78 \pm 1.71
SDCN	0.267 \pm 0.25	1.53 \pm 0.25	1.25 \pm 0.83	1.27 \pm 0.52	4.51 \pm 1.15	7.43 \pm 2.14
DMoN	0.271 \pm 0.45	1.58 \pm 0.95	1.27 \pm 0.93	1.25 \pm 0.82	4.51 \pm 1.12	7.58 \pm 1.35

demonstrate that the dominant scaling behaviour is linear with only minor non-linear components that do not substantially impact the overall scalability characteristics of the model.

In Table 6, we show the time requires for each method for each dataset. It is evident that temporal graph methods are more expensive as these methods model graph sequences whereas a static methods consider the graph structure at the final timestamp with no temporal module. Our proposed approach TGRAIL consistently almost faster than two stage clustering techniques such as TGN across datasets. The reason is clear- TGN requires temporal embedding computation, memory update and perform a post-hoc clustering. Compared to this, TGRAIL performs direct clustering inference in a single model avoiding additional clustering stage after learning the node representations.