

Reducing Misclassification Risk in Dynamic Graph Neural Networks through Abstention

Jayadratha Gayen^{1,2}, Himanshu Pal^{1,2}, Naresh Manwani^{1,3}, and Charu Sharma^{1,3}

¹ Machine Learning Lab, IIIT Hyderabad, India
<https://mll.iiit.ac.in/>

² {jayadratha.gayen,himanshu.pal}@research.iiit.ac.in

³ {naresh.manwani,charu.sharma}@iiit.ac.in

Abstract. Many real-world systems can be modeled as dynamic graphs, where nodes and edges evolve over time, requiring specialized models to capture their evolving dynamics in risk-sensitive applications effectively. Graph neural networks (GNNs) for temporal graphs are one such category of specialized models. For the first time, our approach integrates a reject option strategy within the framework of GNNs for continuous-time dynamic graphs (CTDGs). This allows the model to strategically abstain from making predictions when the uncertainty is high and confidence is low, thus minimizing the risk of critical misclassification and enhancing the results and reliability. We propose a coverage-based abstention prediction model to implement the reject option that maximizes prediction within a specified coverage. It improves the prediction score for link prediction and node classification tasks. Temporal GNNs deal with extremely skewed datasets for the next state prediction or node classification task. In the case of class imbalance, our method can be further tuned to provide a higher weight to the minority class. Exhaustive experiments are presented on four datasets for dynamic link prediction and two datasets for dynamic node classification tasks. This demonstrates the effectiveness of our approach in improving the reliability and area under the curve (AUC)/average precision (AP) scores for predictions in dynamic graph scenarios. The results highlight our model’s ability to efficiently handle the trade-offs between prediction confidence and coverage, making it a dependable solution for applications requiring high precision in dynamic and uncertain environments. Our code is available at: https://github.com/Jayadratha/Cover_DyG.

Keywords: Graph Neural Networks, Temporal Graphs, Reject Option Classification, Link Prediction, Node Classification.

1 Introduction

In the modern era, numerous systems are modeled as dynamic graphs where nodes and edges evolve over time. These systems include social and interaction networks [18], traffic networks [1], trade networks [23], biological networks [2] and

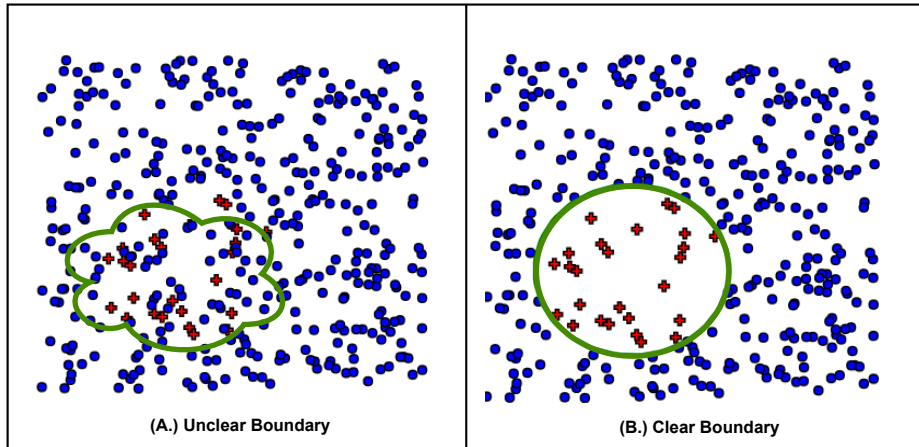


Fig. 1: Decision Boundary: (A) is not clear and smooth; (B) becomes clearer after abstaining confusing/noisy samples. The examples to be rejected are decided by the objective function, which consists of the losses over unrejected samples and a fixed cost for rejected samples. The minimum value of the objective will be achieved only with a proper choice of rejected examples.

transaction networks [20] among others. Particularly in risk-sensitive applications such as fraud detection [24], fake news, polarization identification [6], financial transactions [20], blockchain security [3], epidemic modeling [30], anomaly detection [25] and disease prediction [3], the stakes for accurate and reliable predictions are exceptionally high.

Traditional graph neural networks (GNNs) and their extensions to dynamic graphs have shown promise in capturing the complex interactions and evolving structures within these networks. However, they often fall short where the cost of misclassification is significant and abstaining from making a prediction could mitigate risk. The additional challenge of highly skewed temporal⁴ datasets for node classification (such as Wikipedia, Reddit datasets [18]) makes it harder for the models to predict correctly. Recent works on temporal GNNs [7, 26, 32, 33] solve for node classification and link prediction tasks.

Real-world classification tasks, such as anomaly detection or credit card fraud prediction in dynamic networks, often involve overlapping class distributions and class imbalance, leading to prediction uncertainty. Figure 1(A) illustrates such a scenario. The minority class (red crosses) is infrequent and significantly overlaps with the majority class (blue circles) in a specific region. A classifier attempting to perfectly separate these classes might require a complex decision boundary (conceptualized by the irregular green line), potentially leading to overfitting or unreliable predictions for samples near this boundary.

⁴ “temporal” and “dynamic” will be used interchangeably.

Figure 1(B) conceptually motivates the benefit of abstention. By identifying samples with high prediction uncertainty (typically those within the overlap region), the model can choose to abstain from making a prediction for them. This allows the model to focus on making high-confidence predictions for the remaining samples, which might now be effectively separated by a simpler, more robust decision boundary (conceptualized by the circular green line). This improves both interpretability and trustworthiness, especially in settings where wrong decisions are costly.

Classifiers with abstention are well explored in the machine learning community [4, 12, 15, 19]. Such models allow a classifier to withhold a prediction when the uncertainty associated with the decision is high. This approach has not been widely investigated in the context of dynamic graphs, where the temporal evolution of the data adds a layer of complexity to uncertainty management.

Motivated by the critical need for precision in predictions, especially in environments where trade-offs between false positives and false negatives carry substantial consequences, we explore the integration of the abstention option with the continuous time dynamic graphs (CTDGs) learning method. The open problem, therefore, is to incorporate a mechanism that allows CTDGs to abstain from making predictions under high uncertainty, thus reducing the risk of misclassifications in sensitive applications. To address this, we propose a novel framework that integrates the reject option into CTDGs for both link prediction and node classification tasks. This framework aims to enhance model reliability by minimizing the risk associated with uncertain predictions.

We solve this challenge by extending the methodologies of temporal GNNs to incorporate a reject option classification. Our approach involves designing a specialized neural network architecture that dynamically adjusts its confidence threshold based on the evolving graph structure and node interactions. This method allows the model to strategically abstain from making predictions when uncertainty is high, leveraging a coverage-based model to optimize the trade-off between accurate prediction and the cost of abstention.

Our main goal is uncertainty management using the reject option. While analyzing the data, we find that extreme class imbalance is often present in node classification tasks for CTDG datasets. As uncertainty amplifies due to class imbalance and noisy boundary between the classes, our approach tries to handle both problems together. Previous works have struggled to effectively manage the disproportion between classes, which is particularly problematic in risk-sensitive domains where minority classes may carry significant importance. Our approach handles extreme class imbalance by optimizing the model’s training to better represent minority classes, thereby significantly improving performance in these challenging scenarios. This advancement is a pioneering step in dynamic GNNs research, offering a powerful tool for domains where class imbalance profoundly affects the utility and reliability of the model. Our main contributions are as follows:

1. **Uncertainty management with reject option:** We integrate a reject option strategy into the framework of CTDGs. This empowers the model

to abstain from making predictions when it lacks sufficient confidence, addressing the critical need for risk mitigation in sensitive applications. To achieve this, we propose a coverage-based approach to incorporate the reject option, allowing for the customization of the model’s behavior to optimize predictions within a set coverage limit.

2. **Managing extreme class imbalance:** Drawing from recent literature for class imbalance mitigation, we address this problem in dynamic graph node classification, an area largely unexplored for extreme imbalance.
3. **Comprehensive evaluation:** We conduct extensive experiments on various dynamic graph datasets for link prediction and node classification to demonstrate the effectiveness of our approach. Our results highlight the ability of our models to effectively manage trade-offs between prediction confidence and coverage while significantly improving performance metrics.

2 Related Work

Temporal Graphs. Dynamic graphs [16] are categorized into discrete-time (DTDGs) [22] and continuous-time (CTDGs) [18] settings. In this paper, we focus on CTDGs. CTDGs are addressed using node-based methods, such as TGN [26], and TCL [31] utilize node information, such as temporal neighbors and previous histories of nodes, to create node embedding, and edge-based methods like GraphMixer [7] and CAWN [32] directly generate embeddings for the edge of interest. DyGFormer [34] proposed a transformer-based architecture that uses a neighbor co-occurrence encoding scheme along with a patching technique to effectively capture long-term temporal dependencies in dynamic graphs. Recent methods, such as DyGMamba [9] and FreeDyG [28], incorporate state-space models and frequency-domain analysis for better encoding of interactions. TGB [14] benchmarks popular temporal graph models treated the task as a ranking problem. DyGLib [31] based on PyG [11], provides a unified library for implementing various CTDG methods. We utilize DyGLib to set up experiments and integrate a rejection module.

Uncertainty Estimation & Classification With Rejection (CwR). CwR can be classified broadly into cost-based and coverage-based methods. Cost-based CwR method aligns well with cost-sensitive learning [10] where the cost of misclassification can be used to inform the decision to abstain from prediction. Variants of support vector machine (SVM) with reject option are presented in [27]. [4] introduces a novel approach to multi-class classification with rejection compatible with arbitrary loss functions, addressing the need for flexibility in adapting to different datasets. In coverage-based CwR, SelectiveNet [12] introduces a novel approach by integrating the reject option directly within the deep neural network architecture. [5, 21] provide calibrated losses for multi-class reject option classification. [4] proposes a general recipe to convert any multi-class loss function to accommodate the reject option, calibrated to loss l_{0d1} . They treat rejection as another class. Conformal prediction (CP) provides valid uncertainty estimates by constructing prediction sets that contain the true label

with a predefined confidence level. CF-GNN [13] integrates CP with GNNs to generate uncertainty-aware predictions on graph data. In our task for CTDGs, we only handle binary classification, unlike [8]. So, there is little scope to apply this method in this setting.

3 Preliminaries

This section introduces the concepts and notations used throughout this paper. Our work focuses on extending GNN to temporal graphs, incorporating a classification with a rejection option for enhanced prediction in both link prediction and node classification tasks.

Deep Learning on Dynamic Graphs. Dynamic graphs are a sequence of non-decreasing chronological interactions to accommodate changes over time. They represent evolving relationships within the graph. A dynamic graph is represented as a series of graphs $\mathcal{G}_t = (\mathcal{V}_t, \mathcal{E}_t)$ at discrete time steps t , or as a continuous stream of interactions $(u, v, t) \in \mathcal{E}_t$, where $u, v \in \mathcal{V}_t$ are nodes, and t represents the time of interaction. The primary goal in deep learning on dynamic graphs is to learn a function $f : \mathcal{V}_t \rightarrow \mathbb{R}^d$ that captures not only the structural features but also its temporal dynamics, facilitating tasks such as temporal link prediction and dynamic node classification.

Classification with Rejection. Rejection classification introduces a decision framework in which a model, given an input $x \in \mathcal{X}$, can choose to abstain from making a prediction if it lacks confidence. Given a prediction model $f : \mathcal{X} \rightarrow \mathcal{Y}$, where \mathcal{X} is the input space and \mathcal{Y} is the output label space, and an abstention function $q : \mathcal{X} \rightarrow \{0, 1\}$. For each sample x , the model predicts an abstention score $a(x)$ where $a : \mathcal{X} \in [0, 1]$. We first sort the abstention scores for all the interactions and then select a threshold θ value that aligns with the desired coverage. For an input x , if the abstention score, $a(x) \leq \theta$ then $q(x) = 0$ (the model is confident) and the model predicts $f(x)$. In case $a(x) > \theta$ then $q(x) = 1$ and the model abstains the decision for input x . For a given $\theta \in (0, 1)$, the objective is to maximize the probability of correct prediction for the subset of the inputs where $q(x) = 0$:

$$\max_{f, q} P(f(x) = y | q(x) = 0) = \max_{f, q} P(f(x) = y | a(x) \leq \theta) \quad (1)$$

Problem Formulation. Given a continuous-time temporal graph \mathcal{G}_t with a tuple $(u, v, t) \in \mathcal{E}_t$ representing an interaction between nodes $u, v \in \mathcal{V}_t$, at time $t \in \mathcal{T}$, our objective is to design a temporal GNN model that addresses the problem of learning effective representations for link prediction and node classification while incorporating a reject option to manage uncertainty in predictions.

1. Link Prediction: Given a temporal graph, i.e., source node, destination node, current timestamp, and historical interactions before t , predict the presence or absence of links at time t . Where the abstain mechanism aims to make a prediction or refuse to predict for a pair of nodes based on a coverage criterion.

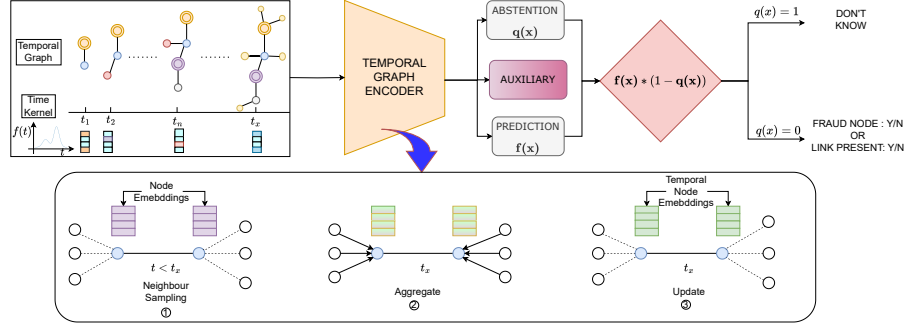


Fig. 2: The model overview of classification with rejection for task-agnostic temporal graphs. Input to the model is a temporal graph with t timestamps, which is processed through an encoder that samples neighbours from all interactions from previous timestamps $t < t_x$ of these nodes, then aggregates and yields an updated temporal node embedding. These embeddings are passed to a *Prediction* function $f(x)$ to predict downstream tasks such as whether a link exists or not for link prediction or whether a user is blocked or not for node classification at time t . The output embeddings also pass through the *Abstention* function, predict a single score, and check if $a(x) > \theta$ threshold i.e. $q(x) = 1$. When it is true, the model abstains; otherwise, it proceeds to make a prediction. The *Auxiliary* head exposes all samples to the model, including samples with high abstention scores, only during training.

2. Node Classification: Predict the label or state of a node in a temporal graph at a time t based on previous information, with the option to abstain from the prediction on certain nodes to ensure high confidence in the predictions made.

4 Method

In this section, we describe our approach for both link prediction and node classification in dynamic graphs in continuous time using the coverage-based method. A neural network architecture designed for dynamic graphs can be seen as a combination of an encoder and a decoder. The encoder serves to transform a dynamic graph into node embeddings, whereas the decoder utilizes these embeddings to make predictions. Our method considers temporal graph networks to incorporate a reject option, enhancing prediction and handling uncertainty effectively. The architecture is illustrated in Fig. 2.

4.1 Temporal Graph Encoder

Our approach is based on the concept of temporal graph networks, which are designed as an encoder-decoder pair for dynamic graph learning. The encoder maps a dynamic graph into the temporal node embeddings $\mathbf{z}_u(t)$, capturing the

temporal interactions among the nodes. The encoder updates the node embeddings with each interaction, ensuring those reflect the graph’s latest state. Few such models are memory-based TGN [26], MLP-based GraphMixer [7], graph transformer-based TCL [31], random-walk-based CAWN [32], SSM-based DyG-Mamba [9] etc.

4.2 Coverage-Based Dynamic Link Prediction

Our approach to link prediction uses the concept of abstention prediction [12] to manage the uncertainty inherent in predicting links in a dynamic graph employing an abstention model (f, q) where f is the prediction function and q is the abstention function. The abstention function $q : \mathcal{X} \rightarrow \{0, 1\}$, if it takes value one, then the prediction is kept on hold for a given pair of nodes based on a coverage criterion. We modify the decoder part of the model to incorporate a reject option through an abstention prediction mechanism. This mechanism evaluates each potential link for its likelihood and the model’s confidence in its prediction. By allowing the model to abstain from making predictions when the uncertainty exceeds a threshold θ , we ensure it only predicts links it is confident about. The abstention prediction objective is defined as follows:

$$\mathcal{L}_{(f,q)}^{\mathcal{E}_t} = \hat{r}(f, q|\mathcal{E}_t) + \lambda \Psi(c - \hat{\phi}(q|\mathcal{E}_t)) \quad (2)$$

where $\hat{r}(f, q|\mathcal{E}_t)$ is the empirical abstention risk, λ is a hyperparameter controlling the importance of the coverage constraint, $\Psi(b) = \max(0, b)^2$ is a quadratic penalty function that helps ensure that the model does not abstain excessively, and c is the target coverage rate. The empirical abstention risk $(\hat{r}(f, q|\mathcal{E}_t))$ and the empirical coverage $(\hat{\phi}(q|\mathcal{E}_t))$ are calculated as follows:

$$\hat{r}(f, q|\mathcal{E}_t) = \frac{\sum_{(u,v,t) \in \mathcal{E}_t} \ell(f(\mathbf{z}_u(t), \mathbf{z}_v(t)), y_{\mathcal{E}_t})(1 - q(\mathbf{z}_u(t), \mathbf{z}_v(t)))}{|\mathcal{E}_t| \hat{\phi}(q|\mathcal{E}_t)} \quad (3)$$

$$\hat{\phi}(q|\mathcal{E}_t) = \frac{1}{|\mathcal{E}_t|} \sum_{(u,v,t) \in \mathcal{E}_t} (1 - q(\mathbf{z}_u(t), \mathbf{z}_v(t))) \quad (4)$$

where ℓ is the *loss function* and $\mathbf{z}_u(t)$, $\mathbf{z}_v(t)$ are the temporal node embeddings for source and destination node for link prediction, respectively, and $y_{\mathcal{E}_t}$ is the binary label of link between two nodes if exists or not. $\hat{\phi}$, the empirical coverage rate of the abstention function q , calculates the rate at which the edges till time t are selected, where $|\mathcal{E}_t|$ is the number of edges. The overall training objective combines the abstention loss with an auxiliary loss. For our case, we use the standard binary cross-entropy as an auxiliary loss, denoted as $\mathcal{L}_h^{\mathcal{E}_t}$ computed on the same prediction task. This auxiliary loss is generated by a dedicated auxiliary prediction head, active only during training. Its purpose is to ensure the shared encoder learns from all training instances, including those the abstention head might abstain from. This approach prevents overfitting to only

high-confidence samples and promotes more robust representations. The overall training objective is:

$$\mathcal{L}^{\mathcal{E}_t} = \alpha \mathcal{L}_{(f,q)}^{\mathcal{E}_t} + (1 - \alpha) \mathcal{L}_h^{\mathcal{E}_t} \quad (5)$$

Where α is a parameter controlling the trade-off between abstention and auxiliary losses.

The temporal graph encoder captures node interaction timing using a time encoder, which is incorporated into the temporal node embeddings. This time information is then passed to both the abstention and prediction functions, allowing the model to decide whether to abstain from making a prediction at a given time t .

4.3 Coverage-Based Dynamic Node Classification

For node classification or user state change prediction [18], we apply a framework similar to link prediction, using the abstention prediction mechanism to introduce a reject option. Node embedding $\mathbf{z}_u(t)$, derived from the temporal graph encoder, incorporates information from the node’s interaction history and its neighbors to produce a comprehensive representation of its current state in the graph. It serves as input to a classification model that incorporates a coverage-based abstention mechanism. The classifier is designed to accurately predict node labels $y_t \in \mathcal{Y}$ while being able to abstain from predictions when uncertainty exceeds the threshold θ . The abstention prediction objective in this case is as follows:

$$\mathcal{L}_{(f,q)}^{\mathcal{V}_t} = \hat{r}(f, q|\mathcal{V}_t) + \lambda \Psi(c - \hat{\phi}(q|\mathcal{V}_t)) \quad (6)$$

Where the empirical abstention risk ($\hat{r}(f, q|\mathcal{V}_t)$) and the empirical coverage ($\hat{\phi}(q|\mathcal{V}_t)$) are calculated as follows:

$$\hat{r}(f, q|\mathcal{V}_t) = \frac{\sum_{u \in \mathcal{V}_t} \ell(f(\mathbf{z}_u(t)), y_t)(1 - q(\mathbf{z}_u(t)))}{|\mathcal{V}_t| \hat{\phi}(q|\mathcal{V}_t)} \quad (7)$$

$$\hat{\phi}(q|\mathcal{V}_t) = \frac{1}{|\mathcal{V}_t|} \sum_{u \in \mathcal{V}_t} (1 - q(\mathbf{z}_u(t))) \quad (8)$$

where $|\mathcal{V}_t|$ is the number of nodes. The node classification model is trained by minimizing a loss function that balances the precision of the abstention prediction with the desired coverage level, which is given by:

$$\mathcal{L}^{\mathcal{V}_t} = \alpha \mathcal{L}_{(f,q)}^{\mathcal{V}_t} + (1 - \alpha) \mathcal{L}_h^{\mathcal{V}_t} \quad (9)$$

where $\mathcal{L}_h^{\mathcal{V}_t}$ is the auxiliary loss (binary cross-entropy).

4.4 Handling Extreme Class Imbalance for Dynamic Node Classification

Dealing with an extreme class imbalance in node classification within CTDGs presents significant challenges. Class imbalance adversely affects model performance, particularly for minority classes, which is crucial for many real-world

applications. To address this, we introduce an approach that modifies the training process to better represent minority classes, thereby improving the model’s performance. We modify the components of our training objective given in Equation (9).

Auxiliary Loss for Class Imbalance: We propose an auxiliary loss function that explicitly accounts for the disproportionate representation of classes by assigning a higher weight to the minority class, which is defined as:

$$\mathcal{L}_h^{\mathcal{V}_t} = \beta \mathcal{L}_{h_{minor}}^{\mathcal{V}_t} + \mathcal{L}_{h_{major}}^{\mathcal{V}_t} \quad (10)$$

where $\mathcal{L}_{h_{minor}}^{\mathcal{V}_t}$ and $\mathcal{L}_{h_{major}}^{\mathcal{V}_t}$ represent the loss for the minority and majority classes, respectively, and β is a weighing factor that amplifies the contribution of the minority class to the overall loss. In our experiments, we set hyperparameter $\beta > 1$, highlighting our focus on the minority class, given its significantly lower representation (less than 0.2% in our dataset). This allows us to tune our model’s performance, ensuring that it provides reliable predictions in the case of extreme class imbalance.

5 Experimental Setup

In this section, we discuss experimental setup, which comprises baselines, datasets used, evaluation tasks performed, metrics used, and implementation details.

5.1 Dataset Used

We use four distinct datasets for dynamic link prediction and two for dynamic node classification where labels are available. We have chosen popular datasets like Wikipedia [18], Reddit (Soical Media) [18], and also included challenging datasets like UN Trade (Economics) and Can. Parl. (Politics) for diversity. These datasets are publicly available ⁵ by DGB [23].

5.2 Baselines Used

Although our framework is compatible with various temporal graph models supported by DyGLib [34], we opted to utilize TGN [26] and GraphMixer [7] due to their widespread adoption as baselines.

TGN Encoder [26] : The key components of TGN Encoder are: a) **Memory module** maintains a compressed representation of each node’s historical interactions. b) **Message function** computes the impact of new interactions on the state of the node. c) **Message aggregator** combines messages from multiple events involving the same node in a batch, improving the update process. d) **Memory updater** integrates new interaction messages into the node’s memory. e) **Embedding module** generates the current embedding of a node using its memory and the memories of its neighbors.

⁵ <https://zenodo.org/record/7213796#.Y1cO6y8r30o>

GraphMixer Encoder [7]:

GraphMixer presents a simple, straightforward, yet effective approach built upon two primary modules: a link encoder and a node encoder based on multilayer perceptrons (MLPs) without relying on complex GNN architectures. a) **Link-encoder** employs a fixed time-encoding function $\cos(\omega t)$. It encodes temporal information of links, followed by an MLP-Mixer [29] to summarize this temporal link information effectively. b) **Node-encoder** utilizes neighbor mean-pooling to aggregate and summarize the features of nodes, capturing essential node identity and feature information for subsequent processing steps.

5.3 Evaluation Tasks and Metrics

We evaluate our approach on two key dynamic graph tasks: link prediction and node classification, aligning with prior work [18, 23, 34].

For link prediction, we predict the likelihood of an edge forming between two nodes at time t , using a 2-layer MLP over concatenated node embeddings. We evaluate both transductive (known nodes) and inductive (unseen nodes) settings, with Average Precision (AP) and AUC-ROC as metrics. Following [23], we use random (rnd), historical (hist: edge re-occurrence), and inductive (ind: unseen test edges) negative sampling strategies (NSS).

For node classification, we fine-tune a pretrained encoder and train a separate MLP for label prediction with a reject option. We report AUC-ROC, suitable for the common class imbalance in these tasks. Datasets without dynamic node labels are excluded.

For both tasks, we split all datasets chronologically: 70% train, 15% validation, 15% test.

5.4 Implementation Details

For the link prediction task, we optimize both models by the Adam optimizer [17] for all datasets. For the node classification task, we optimize both models by the Adam optimizer for the Wikipedia dataset and the SGD for the Reddit dataset. We train the models for 75 epochs and use the early stopping strategy with a patience of 10. We select the model that achieves the best performance in the validation set for testing. We set the batch size to 200, λ to 32, α to 0.5, i.e., equal weight to both abstention prediction loss and auxiliary loss for all the methods on all the datasets. We perform the grid search to find the best settings of some critical hyperparameters. This step is conducted during training. We tune the learning rate using cross-validation. Even though the training process aims to meet the target coverage c via the loss function, the actual coverage achieved on the test set with a fixed threshold *e.g.* $\theta = 0.5$ may differ due to distribution shift. Therefore, to analyze the model’s performance-reliability trade-off consistently, we evaluate at predefined coverage levels (e.g., 90%, 80%, etc.). For each target coverage level reported (e.g., 80%), we calculate the corresponding threshold θ_{80} by finding the 80th percentile of the uncertainty scores $a(x)$ produced by the model on the entire test set. Metrics are then computed solely on the test samples

satisfying $a(x) \leq \theta_{80}$. The θ values range from 0 to 1. For node classification, we search for the best value of β in the range of [2,100] and get the best results when $\beta = 2$ or 5. Both encoders use 100D time encoding and 172D output representation. TGN uses 172D node memory with 2 graph attention heads, using GRU for memory updates. GraphMixer uses 2 MLP-Mixer layers with a time slot gap of 2000 interactions. We run each model five times with seeds from 0 to 4 and report the average performance and standard deviation to eliminate deviations. We use NVIDIA GeForce RTX 4090 with 24 GB memory for link prediction task. The GPU device used for the node classification task is NVIDIA GeForce RTX 2080 Ti with 11 GB of memory.

6 Experimental Results

In this section, we report the performance of our proposed continuous-time dynamic graph neural network with a reject option and compare our results with baselines.

Table 1: AP for transductive dynamic link prediction with random, historical, and inductive negative sampling strategies for various coverage rates. The best and second-best results are emphasized by **bold** and underlined.

NSS	Coverage (%)	TGN				GraphMixer			
		Wikipedia	Reddit	UN Trade	Can. Parl.	Wikipedia	Reddit	UN Trade	Can. Parl.
rnd	100	98.56 \pm 0.06	98.63 \pm 0.02	64.87 \pm 1.93	74.14 \pm 1.51	97.30 \pm 0.05	97.35 \pm 0.02	61.68 \pm 0.16	79.96 \pm 0.95
	90	99.33 \pm 0.06	99.28 \pm 0.07	71.67 \pm 1.04	75.38 \pm 1.44	98.24 \pm 0.07	97.66 \pm 0.06	68.57 \pm 0.21	80.56 \pm 1.71
	80	99.73 \pm 0.08	99.54 \pm 0.06	76.59 \pm 1.50	77.21 \pm 5.87	98.89 \pm 0.06	98.41 \pm 0.03	<u>72.50 \pm 1.24</u>	82.44 \pm 1.31
	70	99.85 \pm 0.02	99.81 \pm 0.04	79.09 \pm 1.66	81.57 \pm 7.38	99.61 \pm 0.14	99.32 \pm 0.08	71.64 \pm 8.98	95.25 \pm 0.65
	60	99.88 \pm 0.02	<u>99.86 \pm 0.04</u>	<u>81.46 \pm 1.81</u>	90.19 \pm 3.50	99.68 \pm 0.11	<u>99.52 \pm 0.11</u>	69.99 \pm 11.27	95.98 \pm 0.57
	50	<u>99.87 \pm 0.05</u>	99.91 \pm 0.03	84.46 \pm 0.95	<u>89.46 \pm 6.67</u>	<u>99.63 \pm 0.08</u>	99.65 \pm 0.11	80.94 \pm 2.11	94.68 \pm 1.25
hist	100	87.07 \pm 0.81	80.63 \pm 0.30	59.44 \pm 2.22	70.92 \pm 2.11	91.14 \pm 0.16	77.50 \pm 0.54	57.61 \pm 1.11	80.66 \pm 1.37
	90	90.81 \pm 0.59	81.68 \pm 1.07	70.36 \pm 0.80	74.18 \pm 1.94	94.25 \pm 0.43	78.17 \pm 0.25	63.36 \pm 3.08	83.10 \pm 1.57
	80	94.46 \pm 0.64	83.78 \pm 0.33	74.44 \pm 3.55	72.82 \pm 6.51	96.57 \pm 0.63	80.57 \pm 0.33	69.35 \pm 3.68	84.08 \pm 2.60
	70	95.99 \pm 3.62	87.25 \pm 1.76	71.94 \pm 8.63	76.08 \pm 9.00	<u>97.99 \pm 2.64</u>	83.61 \pm 1.96	75.01 \pm 4.15	<u>94.97 \pm 0.55</u>
	60	<u>97.06 \pm 3.84</u>	<u>88.66 \pm 1.68</u>	<u>76.36 \pm 2.16</u>	<u>86.97 \pm 4.56</u>	98.56 \pm 2.81	<u>85.32 \pm 1.76</u>	69.62 \pm 1.57	95.05 \pm 0.97
	50	97.82 \pm 4.50	89.27 \pm 2.41	76.65 \pm 2.85	87.34 \pm 3.79	97.40 \pm 3.50	85.91 \pm 2.04	81.70 \pm 2.52	94.18 \pm 1.33
ind	100	86.57 \pm 0.96	88.02 \pm 0.35	61.70 \pm 2.21	68.05 \pm 2.69	88.83 \pm 0.13	85.21 \pm 0.28	60.89 \pm 1.01	77.74 \pm 1.42
	90	89.75 \pm 1.03	90.17 \pm 0.82	72.31 \pm 0.77	71.28 \pm 2.90	92.11 \pm 0.23	85.76 \pm 0.34	67.57 \pm 2.52	79.98 \pm 2.47
	80	93.19 \pm 0.57	92.89 \pm 0.80	77.43 \pm 2.94	69.06 \pm 9.77	95.04 \pm 0.37	88.74 \pm 0.43	74.08 \pm 1.72	79.26 \pm 3.90
	70	95.85 \pm 3.61	94.95 \pm 2.27	75.75 \pm 6.69	72.87 \pm 10.05	96.45 \pm 3.70	91.38 \pm 0.86	79.47 \pm 2.22	93.47 \pm 0.56
	60	96.76 \pm 4.00	96.57 \pm 2.68	78.65 \pm 1.93	84.91 \pm 5.02	97.66 \pm 3.90	93.11 \pm 1.45	77.98 \pm 1.54	94.11 \pm 0.84
	50	97.75 \pm 4.63	<u>95.93 \pm 3.33</u>	<u>78.30 \pm 2.15</u>	85.12 \pm 6.52	94.79 \pm 6.29	94.35 \pm 1.60	85.47 \pm 2.57	94.28 \pm 1.07

6.1 Results: Dynamic Link Prediction

We evaluated our rejection-aware framework on the dynamic link prediction task using both TGN and GraphMixer as base encoders, comparing performance against standard baselines run using DyGLib [42]. The primary goal was to assess the improvement in prediction reliability when allowing the model to abstain based on uncertainty.

Table 1 summarizes the transductive link prediction results, reporting Average Precision (AP) across all three negative sampling strategies (NSS). The 100% coverage rows represent the baseline performance without rejection. A consistent and significant trend is observed across all datasets, models, and NSS

settings: as the coverage level decreases (meaning the model abstains on a larger fraction of uncertain predictions, from 90% down to 50% coverage), the AP score calculated on the remaining non-abstained predictions progressively improves.

This shows the effectiveness of the rejection mechanism in filtering out low-confidence predictions, thereby enhancing the precision of the accepted ones. For example, focusing on the challenging Canadian Parliament (Can. Parl.) dataset, which exhibits lower baseline scores, the improvement is notable. With our model, the AP/AUC score notably increases, emphasizing the model’s efficiency in managing harder prediction environments. For example, using TGN with random NSS, the AP increases from a baseline of (1st row) at 100% coverage to 90.19% AP (5th row) at 60% coverage in Can. Parl. dataset, marking almost 16.05% improvement in AP. These gains show the value of abstention for boosting trustworthiness, particularly in environments with sparse or noisy interaction data typical of many real-world social or information networks. Further

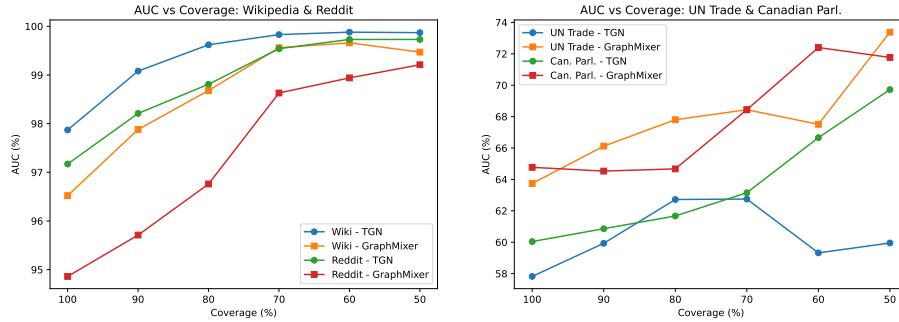


Fig 3: Performance (AUC %) vs. Coverage (%) Trade-off for Inductive Link Prediction in Random NSS.

illustrating the performance-reliability trade-off, Figure 3 plots the relationship between coverage and prediction performance (AUC) for the *inductive* link prediction task (using Random NSS). Consistent with the AP results, the inductive AUC generally increases as coverage decreases (moving rightward on the x-axis). This confirms that the benefits of abstention—achieving higher accuracy on a more reliable subset of predictions—extend to the demanding inductive setting, where generalization to unseen graph elements is crucial. The visual trend highlights the model’s capability to manage uncertainty effectively across different evaluation scenarios.

One observation is that AP/AUC increases to a certain point (peak) as coverage decreases. It starts decreasing as coverage goes on decreasing.

6.2 Results: Dynamic Node Classification

For dynamic node classification, we evaluate our model on the Wikipedia and Reddit datasets, as shown in Table 2. The initial results with 100% coverage

use DyGLib, which establishes a benchmark. As we implement the reject option, progressively rejecting 10% of samples in each subsequent test until we reach 60% coverage, we observe a significant improvement in AUC scores (where $\beta = 1$). For example, in Table 2 we achieve 89.78% AUC (ninth row) in 60% coverage w.r.t. 86.23% AUC (first row) in 100% coverage for the TGN model on Wikipedia, marking almost 3.55% improvement in AUC.

This increase in AUC with decreasing coverage indicates that the model is effectively prioritizing more reliable predictions, as less certain predictions are rejected. Furthermore, due to the extreme imbalance in class distribution (minority class below 0.2%), we provide higher weightage to minority class in the auxiliary loss, which led to further improvements in AUC (where $\beta > 1$). For example, in Table 2, we get 69.58% AUC (tenth row) with $\beta = 5$ w.r.t. 66.03% AUC (ninth row) in 60% coverage for TGN model on Reddit dataset. This marks almost 3.55% improvement in AUC while addressing the class imbalances. This strategy improves the efficacy of our approach in not only managing the reject option but also addressing extreme class imbalances, which are particularly challenging in dynamic environments.

Table 2: AUC for coverage-based dynamic node classification for various coverage rates with and without handling class imbalance. The best and second-best performing results are emphasized by **bold** and underlined fonts. * denotes β is 2, otherwise 5 for all other cases where $\beta > 1$.

Coverage(%)	β	TGN		GraphMixer	
		Wikipedia	Reddit	Wikipedia	Reddit
100	= 1	86.23 \pm 3.30	63.08 \pm 1.48	86.19 \pm 1.10	64.81 \pm 2.08
90	= 1	88.00 \pm 2.00	64.81 \pm 1.54	87.27 \pm 1.72	65.26 \pm 2.07
	> 1	88.44 \pm 2.33*	65.12 \pm 1.53	87.61 \pm 1.06*	66.20 \pm 1.88
80	= 1	88.78 \pm 3.46	65.63 \pm 1.31	87.93 \pm 1.53	66.08 \pm 2.99
	> 1	<u>90.02 \pm 2.29*</u>	66.31 \pm 3.26	<u>88.05 \pm 1.72</u>	67.00 \pm 2.40
70	= 1	89.48 \pm 2.53	66.10 \pm 2.34	87.79 \pm 1.37	67.42 \pm 3.27
	> 1	90.64 \pm 3.37	<u>67.90 \pm 1.81</u>	88.38 \pm 0.99*	68.23 \pm 2.96
60	= 1	89.78 \pm 2.76	66.03 \pm 3.45	85.97 \pm 2.62	67.23 \pm 3.36
	> 1	89.86 \pm 2.52	69.58 \pm 2.96	84.83 \pm 2.30*	<u>67.87 \pm 4.01</u>

In real-world scenarios, we can't reject 40-50% predictions. But as we can see, with a decrease in coverage, performance keeps on increasing. Note that the model rejects examples; it does not predict very confidently. We have to find a balance between the rate of abstention and the model confidence based on the application domain, number of data points, etc. Overall, the experiments validate our hypothesis that integrating a reject option within CTDGs significantly enhances the reliability and accuracy of predictions in dynamic graphs. The clear improvement in performance metrics across different datasets, coverage settings, and tasks confirms the practical utility of our model in risk-sensitive applications, where precision and reliability are crucial.

7 Conclusion

The experimental findings convincingly demonstrate the efficacy of our proposed CTDG model with a reject option for both link prediction and node classification tasks. The model strategically abstains from making predictions when encountering high uncertainty, leading to a significant improvement in AUC/AP scores. This establishes the model’s proficiency in managing the trade-off between prediction confidence and coverage. In the future, we can explore cost-based or other abstention methods to enhance the model’s applicability in different contexts. Another future direction could involve demonstrating the model’s performance in real-world applications, showcasing its practical relevance. Considering our method for another kind of graph model, such as discrete time temporal graphs, can also be explored. Overall, our work presents a compelling solution for dynamic graph applications that demand high precision and reliability, particularly in risk-sensitive domains. The model’s ability to effectively manage uncertainty and class imbalance makes it a valuable tool for various real-world applications.

References

1. Bai, L., Yao, L., Li, C., Wang, X., Wang, C.: Adaptive graph convolutional recurrent network for traffic forecasting. In: NeurIPS (2020)
2. Behrouz, A., Hashemi, F.: Learning temporal higher-order patterns to detect anomalous brain activity. In: TGLW @ NeurIPS (2023)
3. Behrouz, A., Seltzer, M.: Anomaly detection in multiplex dynamic networks: from blockchain security to brain disease prediction. In: TGLW @ NeurIPS (2022)
4. Cao, Y., Cai, T., Feng, L., Gu, L., GU, J., An, B., Niu, G., Sugiyama, M.: Generalizing consistent multi-class classification with rejection to be compatible with arbitrary losses. In: NeurIPS (2022)
5. Charoenphakdee, N., Cui, Z., Zhang, Y., Sugiyama, M.: Classification with rejection based on cost-sensitive classification. In: ICML. pp. 1507–1517 (2021)
6. Chomel, V., Cuvelles-Magar, N., Panahi, M., Chavalarias, D.: Polarization identification on multiple timescale using representation learning on temporal graphs in eulerian description. In: TGLW @ NeurIPS (2022)
7. Cong, W., Zhang, S., Kang, J., Yuan, B., Wu, H., Zhou, X., Tong, H., Mahdavi, M.: Do we really need complicated model architectures for temporal networks? In: ICLR (2023)
8. Davis, E., Gallagher, I., Lawson, D.J., Rubin-Delanchy, P.: Valid conformal prediction for dynamic GNNs. In: ICLR (2025)
9. Ding, Z., Li, Y., He, Y., Norelli, A., Wu, J., Tresp, V., Ma, Y., Bronstein, M.: Dygmamba: Efficiently modeling long-term temporal dependency on continuous-time dynamic graphs with state space models. arXiv preprint arXiv:2408.04713 (2024)
10. El-Yaniv, R., Wiener, R.: On the foundations of cost-sensitive learning. JMLR pp. 1393–1426 (2007)
11. Fey, M., Lenssen, J.E.: Fast graph representation learning with PyTorch Geometric. In: RLGM @ ICLR-W (2019)
12. Geifman, Y., El-Yaniv, R.: SelectiveNet: A deep neural network with an integrated reject option. In: ICML (2019)

13. Huang, K., Jin, Y., Candes, E., Leskovec, J.: Uncertainty quantification over graph with conformalized graph neural networks. *NeurIPS* **36** (2024)
14. Huang, S., Poursafaei, F., Danovitch, J., Fey, M., Hu, W., Rossi, E., Leskovec, J., Bronstein, M.M., Rabusseau, G., Rabbany, R.: Temporal graph benchmark for machine learning on temporal graphs. In: *NeurIPS Datasets and Benchmarks Track* (2023)
15. Kalra, B., Shah, K., Manwani, N.: Risan: Robust instance specific deep abstention network. In: *UAI* (2021)
16. Kazemi, S.M., Goel, R., Jain, K., Kobyzev, I., Sethi, A., Forsyth, P., Poupart, P.: Representation learning for dynamic graphs: A survey. *JMLR* pp. 1–73 (2020)
17. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: *ICLR* (2015)
18. Kumar, S., Zhang, X., Leskovec, J.: Predicting dynamic embedding trajectory in temporal interaction networks. In: *KDD*. pp. 1269–1278 (2019)
19. Manwani, N., Desai, K., Sasidharan, S., Sundararajan, R.: Double ramp loss based reject option classifier. In: *PAKDD*. pp. 151–163 (2015)
20. Nath, P., Waghmare, G., Agrawal, N., Kumar, N., Asthana, S.: TBoost: Gradient boosting temporal graph neural networks. In: *TGLW @ NeurIPS* (2023)
21. Ni, C., Charoenphakdee, N., Honda, J., Sugiyama, M.: On the calibration of multiclass classification with rejection. In: *NeurIPS*. vol. 32 (2019)
22. Pareja, A., Domeniconi, G., Chen, J., Ma, T., Suzumura, T., Kanezashi, H., Kaler, T., Schardl, T., Leiserson, C.: Evolvegc: Evolving graph convolutional networks for dynamic graphs. In: *AAAI*. pp. 5363–5370 (2020)
23. Poursafaei, F., Huang, A., Pelrine, K., Rabbany, R.: Towards better evaluation for dynamic link prediction. In: *NeurIPS Datasets and Benchmarks Track* (2022)
24. Reddy, S., Poduval, P., Chauhan, A.V.S., Singh, M., Verma, S., Singh, K., Bhowmik, T.: Tegra: temporal and graph based fraudulent transaction detection framework. In: *ACM International Conference on AI in Finance*. pp. 1–8 (2021)
25. Reha, J., Lovisotto, G., Russo, M., Gravina, A., Grohnfeldt, C.: Anomaly detection in continuous-time temporal provenance graphs. In: *TGLW @ NeurIPS* (2023)
26. Rossi, E., Chamberlain, B., Frasca, F., Eynard, D., Monti, F., Bronstein, M.: Temporal graph networks for deep learning on dynamic graphs. In: *GRL @ ICML-W* (2020)
27. Shah, K., Manwani, N.: Sparse reject option classifier using successive linear programming. In: *AAAI*. pp. 4870–4877 (2019)
28. Tian, Y., Qi, Y., Guo, F.: Freedyg: Frequency enhanced continuous-time dynamic graph model for link prediction. In: *ICLR* (2024)
29. Tolstikhin, I.O., Houlsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A., Keysers, D., Uszkoreit, J., Lucic, M., Dosovitskiy, A.: Mlp-mixer: An all-mlp architecture for vision. In: *NeurIPS*. pp. 24261–24272 (2021)
30. Varugunda, S.S., Fan, C.H., Wang, L.: Exploring graph structure in graph neural networks for epidemic forecasting. In: *TGLW @ NeurIPS* (2023)
31. Wang, L., Chang, X., Li, S., Chu, Y., Li, H., Zhang, W., He, X., Song, L., Zhou, J., Yang, H.: TCL: transformer-based dynamic graph modelling via contrastive learning. *CoRR* (2021)
32. Wang, Y., Chang, Y., Liu, Y., Leskovec, J., Li, P.: Inductive representation learning in temporal networks via causal anonymous walks. In: *ICLR* (2021)
33. Xu, D., Ruan, C., Körpeoglu, E., Kumar, S., Achan, K.: Inductive representation learning on temporal graphs. In: *ICLR* (2020)
34. Yu, L., Sun, L., Du, B., Lv, W.: Towards better dynamic graph learning: New architecture and unified library. In: *NeurIPS* (2023)