

Learning Point Processes using Recurrent Graph Network

Saurabh Dash
saurabhdash@gatech.edu

Xueyuan She
xshe6@gatech.edu

Saibal Mukhopadhyay
saibal@ece.gatech.edu

Abstract—We present a novel Recurrent Graph Network (RGN) approach for predicting discrete marked event sequences by learning the underlying complex stochastic process. Using the framework of Point Processes, we interpret a marked discrete event sequence as the superposition of different sequences each of a unique type. The nodes of the Graph Network use LSTM to incorporate past information whereas a Graph Attention Network (GAT Network) introduces strong inductive biases to capture the interaction between these different types of events. By changing the self-attention mechanism from attending over past events to attending over event types, we obtain a reduction in time and space complexity from $\mathcal{O}(N^2)$ (total number of events) to $\mathcal{O}(|\mathcal{Y}|^2)$ (number of event types). Experiments show that the proposed approach improves performance in log-likelihood, prediction and goodness-of-fit tasks with lower time and space complexity compared to state-of-the-art Transformer based architectures.

I. INTRODUCTION

Discrete event sequences are ubiquitous around us - hospital visits, tweets, financial transactions, earthquakes, human activity. These sequences are primarily characterized by the event timestamps however, they can also contain markers and other additional information like type and location. Often these sparse observations have complex hidden dynamics responsible for them and when observed through the lens of Point Processes, these sequences can be interpreted as realization of an underlying stochastic process [1]. This underlying point process can be a simple stochastic process like a Poisson process or could be doubly stochastic and have dependencies on both time and history like the Hawkes process [2].

Traditionally these processes have been described using simplified historical dependencies. The Hawkes process assumes additive mutual excitation however, these underlying assumptions are not valid for practical real world data. Recent works seek to leverage the advances in deep learning to relax these restrictions by defining a rich, flexible family of models to incorporate historical information [3], [4], [5], [6], [7]. Recurrent Neural Network (RNN) [3], [4], [7] and Transformer based [5], [6] architectures have been explored to capture historical dependencies to better predict the time and type of future events compared to a pre-specified Hawkes Process. However, shallow RNNs only capture simple state transitions and multi-layer RNNs are difficult to train. On the other hand, Transformer

based approaches are highly expressive and interpretable but the multi-headed self-attention mechanism consumes exorbitant amounts of memory for large sequences.

Our novel approach simultaneously solves the problems faced by the above approaches. An entire discrete event sequence with marks can be interpreted as a superposition of multiple sequences each only containing events of a single unique marker. We process the evolution of each of these unique event type sequences with an LSTM [8]. These different LSTM hidden states are then used to generate a dynamic relational graph and model complex dependencies between event types using a Graph Attention Network [9] as shown in figure 1. A key innovation of our architecture lies in attending over event types instead of all past events. This leads to reduction in space (memory) and time (number of computations) complexity as instead of calculating attention over all past inputs - which can get prohibitively expensive as the sequences get longer $\mathcal{O}(N^2)$, we attend over event types to obtain a time-varying attention matrix with constant time and space complexity in the order of $\mathcal{O}(|\mathcal{Y}|^2)$, independent of the sequence length. Additionally, this time-varying event attention matrix helps with interpretability as it encodes the evolution of the underlying relational graph between event types as new events are observed in the data stream. To the best of our knowledge, this is the first time a Recurrent Graph Network using expressive attention mechanism to encode a dynamic relational graph is used for learning Point Processes. The advantages of the proposed method over previous approaches is summarized in table I.

The major contributions of this paper are summarised as follows:

- We present a Recurrent Graph Network (RGN) approach to learn point processes that uses strong inductive biases to better encode historical information to model the conditional intensity function.
- We show that the proposed model performs better in log-likelihood and predictive tasks for multiple datasets.
- We show that this formulation allows for interpretability of the underlying temporally evolving Relational Graph between event types.
- We show that RGN has lower computational complexity and activation footprint than state-of-the-art Transformer based approaches leading to computa-

Table I: Comparison with previous approaches

Model	Stack Multiple Layers	Small Activation Footprint	Model Interpretability
RMTTP [3]	✗	✓	✗
NHP [4]	✗	✓	✗
THP [5]	✓	✗	✓
This work	✓	✓	✓

tion and memory savings.

We perform experiments on standard point process datasets like - Retweets [10], Financial Transactions [3], StackOverflow [11] and MIMIC-II [12]. In addition to these, we also introduce a new dataset for benchmarking point process models extracted from StarCraft II game replays. The underlying complex stochastic process that generates this data is the high level strategy that the players adopt over the course of the game. The Starcraft II dataset can be found here¹ and the supplementary materials can be found here².

Our experiments show that RGN improves over state-of-the-art Transformer Hawkes Process (THP) on log-likelihood by upto 8%, event type prediction accuracy by upto 4% and event time error by upto 13%.

II. RELATED WORK

Du et al. [3] propose using an RNN (RMTTP) to encode the history to a hidden state \mathbf{h}_t that is used to define the conditional intensity $\lambda^*(t) = \exp(\mathbf{v}^T \mathbf{h}_t + w(t - t_i) + b)$. Mei & Eisner [4] proposed a novel RNN architecture - Continuous-time LSTM (CLSTM) which modifies the LSTM cell to incorporate a decay δ when an event is not observed to a steady state \bar{c} . Since CLSTM work in continuous time, there is no need to explicitly feed the timestamps of the events like RMTTP. However both models only contain shallow RNNs which only capture simple state transitions [13] as multi-layer RNNs are difficult to train [14] due to vanishing / exploding gradients. To alleviate these concerns Zuo et al. [5] and Zhang et al. [6] make use of multi-headed self attention [15] based Transformer architectures to better distill out historical influence. This allows using multi-layer architectures to model complex interactions compared to RNN based architectures, however the $\mathcal{O}(N^2)$ cost of self attention computation leads to prohibitively large memory footprint for longer sequences. Even if a single sequence in a mini-batch is long, the other sequences have to be zero padded to match the longest length for batch processing, leading to a lot of wasted memory and redundant computation. In contrast, Schur et al. [7] draw upon advances in Normalizing Flows [16] and propose a mixture distribution to define a model using the conditional density function instead of the conditional intensity function. This allows for the conditional density function to be multi-modal

allowing for much more complex distributions; however, the historical influence is again encoded with a simple RNN which leads to degraded performance due to lack to rich historical information.

Chang et al. [17] propose a dynamic message passing neural network named TDIG-MPNN which aims to learn continuous-time dynamic embeddings from sequences of interaction data. Although our proposed network also benefits from learning dynamic graph embeddings, the end goal of this work is to learn the underlying conditional intensity function that is responsible for observed data and not the time-varying interaction graph between different entities. Moreover, the datasets under consideration here do not contain any interaction data. The work closest to the proposed approach is ARNPP-GAT [18]. ARNPP-GAT divides the representations of users (marks) into two categories: long-term representation which is modelled using Graph Attention Layer (GAT Layer) and short-term representation which is modelled using a Gated Recurrent Unit (GRU). Although we use similar blocks, the architecture of the model is drastically different. ARNPP-GAT leverages GAT layer to model a time-independent interaction between users (marks) while the historical context is modelled by a GRU. These embeddings are then combined for predicting the time and location of next event. On the other hand, RGN uses GAT layers at every time-step to model a dynamic relational graph between different event types. Moreover, we also use an additional log-likelihood loss term which encourages the model to learn the underlying stochastic process better.

III. BACKGROUND

A. Point Processes

A temporal point process is a stochastic process composed of a time series of events that occur instantaneously in continuous time [1], [19]. The realization of this point process is a set of strictly increasing arrival times $\tau = \{t_1, \dots, t_N\}$. We are interested in a Marked Point Process where each event also has an associated marker $y \in \mathcal{Y}$. In real world data, event probabilities can be affected by past events, thus the probability of an event occurring in an infinitesimally small window $[t, t + dt)$ is called the conditional intensity $\lambda^*(t) := \lambda(t|\mathcal{H}_t)$. \mathcal{H}_t denotes the history $\{(t', y') \in \tau | t' < t\}$. The * symbol represents the dependence on history. Given the conditional intensity function, the conditional density function can be obtained as shown in Eq. 1 below [20].

$$f^*(t) = \lambda^*(t) \exp\left(-\int_{t_n}^t \lambda^*(s) ds\right) \quad (1)$$

B. Recurrent Networks

Recurrent Networks are autoregressive models that can incorporate past information along with the present input for prediction. RNNs and its modern variants LSTM [8] and GRU [21] have shown remarkable success in sequence

¹<https://figshare.com/s/c028296e953788b25599>

²<https://figshare.com/s/7a10a6aa2c0752f2cca3>

modelling tasks like - language modelling [22], [23], video processing [24] and dynamical systems [25]. At each time step, an input \mathbf{i}_t is fed into the model with the past internal state information \mathbf{h}_t . The internal state is updated and used to predict the output \mathbf{o}_t .

C. Graph Networks

Graph Networks (GN) [26] describe a class of functions for relational reasoning over graph structured data. These take in a graph structured data denoted by the 3-tuple (\mathbf{u}, V, E) . \mathbf{u} denotes the global attributes, $V := \{\mathbf{v}_i\}_{i=1}^{N^v}$ denotes the set of node attributes and $E := \{(\mathbf{e}_i, r_k, s_k)\}_{i=1}^{N^e}$ denotes the set of edges of a graph where \mathbf{e}_k is the edge attribute, r_k is the receiver node and s_k is the sender node.

A full GN block enables flexible representations in terms of representation of attributes and in terms of structure of the graph and can be used to describe a wide variety of architectures [9], [15]. The main unit of computation is the GN block which performs graph-to-graph operations. The GN block takes a graph structure as an input, performs operations over the graph and outputs a graph.

A full GN block contains three update functions ϕ and three aggregate functions ρ [26]:

$$\mathbf{e}'_k = \phi^e(\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}, \mathbf{u}) \quad (2)$$

$$\mathbf{v}'_i = \phi^v(\bar{\mathbf{e}}'_i, \mathbf{v}_i, \mathbf{u}) \quad (3)$$

$$\mathbf{u}' = \phi^u(\bar{\mathbf{e}}', \bar{\mathbf{v}}', \mathbf{u}) \quad (4)$$

$$\bar{\mathbf{e}}'_i = \rho^{e \rightarrow v}(E'_i) \quad (5)$$

$$\bar{\mathbf{e}}' = \rho^{e \rightarrow u}(E') \quad (6)$$

$$\bar{\mathbf{v}}' = \rho^{v \rightarrow u}(V') \quad (7)$$

where, $E'_i := \{(\mathbf{e}'_k, r_k, s_k)\}_{r_k=i, k=1}^{N^e}$ is the set of all updated edge attributes that have node i as the receiver, $V' := \{\mathbf{v}'_i\}_{i=1}^{N^v}$ is the set of updated node attributes and $E' := \bigcup_i E'_i = \{(\mathbf{e}_k, r_k, s_k)\}_{k=1}^{N^e}$ is the set of updated edge attributes. The update functions ϕ are shared across the nodes and edges allowing for a reduction in parameter count and encouraging a form of combinatorial generalization [26]. These can be any arbitrary functions or more generally parameterized using a Neural Network architecture like - Multi Layer Perceptron (MLP) or incorporate past information using a Recurrent Neural Network (RNN). On the other hand, ρ are permutation invariant aggregation functions which take a set of inputs and reduce them to a single aggregate element, for example max or mean.

IV. PROPOSED ARCHITECTURE

A. Recurrent Graph Network for Learning Point Processes

An entire discrete event sequence with marks can be interpreted as a superposition of multiple sequences each only containing events of a single unique marker. Using the input sequence, we build a fully connected graph with $|\mathcal{Y}|$ nodes, each containing a latent embedding for a specific type marker which allows us to understand how different

event types influence future events. As we are building a relational graph from non-graph structured inputs, the relational information in the graph is not explicit and needs to be inferred. This relational graph has node attributes and edge weights evolving with time as new events are observed. Thus in a sequence \mathcal{S}_i , with N_T events, we can interpret the relational graph undergoing N_T transitions which can be processed naturally using a Recurrent Graph Network. We are interested in predicting the conditional intensity $\lambda^*(t)$ which is a global property, which makes the Recurrent Graph Network a graph focused Graph Network [26].

At each event (t_i, y_i) in an event sequence \mathcal{S} , the inputs to the model are the previous node attributes $\{\mathbf{v}_{t_{i-1}}^y\}_{y=1}^{|\mathcal{Y}|}$, $\mathbf{v}_{t_{i-1}}^y \in \mathbb{R}^{d^v}$ and the input embedding \mathbf{x}_{t_i} . Once a new event is observed, there are three types of updates that take place:

1) *Input-to-Graph Update:* Each unique-marked sequence is assigned a designated graph node for processing. When an event of a particular type is observed, we want to first update the node attribute corresponding to this marker. This is performed using an LSTM [8]. The LSTM corresponding to node y_i updates its node attribute $\mathbf{v}_{t_i}^{y_i}$ using the previous node attribute $\mathbf{v}_{t_{i-1}}^{y_i}$ and \mathbf{x}_{t_i} using the following equations, while the rest of the node attributes remain unchanged in this update.

$$\mathbf{f}_{t_i}^{y_i} = \sigma(W_f^y \mathbf{x}_{t_i} + U_f^y \mathbf{v}_{t_{i-1}}^{y_i} + \mathbf{b}_f^{y_i}) \quad (8)$$

$$\mathbf{i}_{t_i}^{y_i} = \sigma(W_i^y \mathbf{x}_{t_i} + U_i^y \mathbf{v}_{t_{i-1}}^{y_i} + \mathbf{b}_i^{y_i}) \quad (9)$$

$$\mathbf{o}_{t_i}^{y_i} = \sigma(W_o^y \mathbf{x}_{t_i} + U_o^y \mathbf{v}_{t_{i-1}}^{y_i} + \mathbf{b}_o^{y_i}) \quad (10)$$

$$\mathbf{g}_{t_i}^{y_i} = \tanh(W_g^y \mathbf{x}_{t_i} + U_g^y \mathbf{v}_{t_{i-1}}^{y_i} + \mathbf{b}_g^{y_i}) \quad (11)$$

$$\mathbf{c}_{t_i}^{y_i} = \mathbf{f}_{t_i}^{y_i} \circ \mathbf{c}_{t_{i-1}}^{y_i} + \mathbf{i}_{t_i}^{y_i} \circ \mathbf{g}_{t_i}^{y_i} \quad (12)$$

$$\mathbf{v}_{t_i}^{y_i} = \mathbf{o}_{t_i}^{y_i} \circ \tanh(\mathbf{c}_{t_i}^{y_i}) \quad (13)$$

2) *Graph-to-Graph Update:* Once the node attribute $\mathbf{v}_{t_i}^{y_i}$ corresponding of the observed event type y_i node has been updated, information needs to be propagated to other nodes to update the node attributes. For this purpose, we propose using Graph Attention Network (GAT) [9] as the attention mechanism can be used to assign edge weights not explicitly present in the relational graph using the node attributes. In this section, we express the Graph Attention Layer [9] in terms of Graph Network [26] operations. For notational simplicity, we drop the temporal dependence of node attributes $\mathbf{v}_{t_i}^y$ at current time t_i and represent them with \mathbf{v}_y .

Edge Update: The edge update function ϕ^e only uses the node attributes of the sender node s_k and receiver node r_k and outputs a scalar $a'_k \in \mathbb{R}^+$, vector $\mathbf{b}'_k \in \mathbb{R}^{d^e}$ tuple. NN represents a multi-layer perceptron (MLP).

$$\mathbf{e}'_k = (a'_k, \mathbf{b}'_k) = \phi^e(e_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}, \mathbf{u}) := f^e(\mathbf{v}_{r_k}, \mathbf{v}_{s_k}) \quad (14)$$

$$a'_k = \exp(\text{NN}_{\alpha'}([\text{NN}_{\alpha}(\mathbf{v}_{r_k}), \text{NN}_{\alpha}(\mathbf{v}_{s_k})])) \quad (15)$$

$$\mathbf{b}'_k = \text{NN}_{\beta}(\mathbf{v}_{s_k}) \quad (16)$$

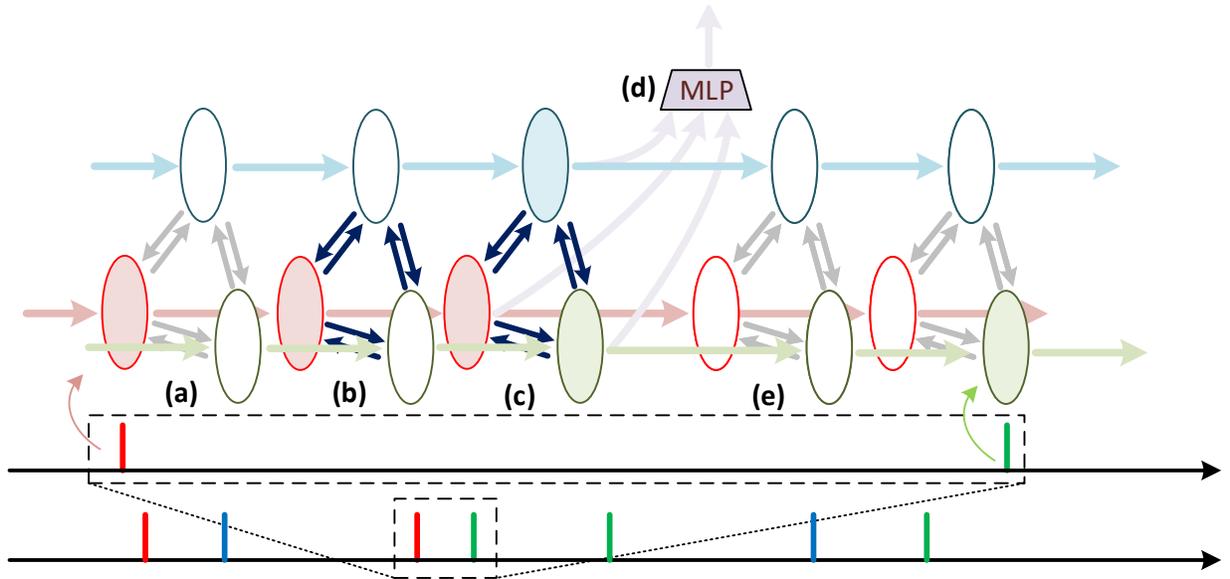


Figure 1: Updates in Recurrent Graph Network: (a) *Input-to-Graph*: The model observes an input event, using the past node attributes, updates event node. (b) *Edge Update*: Based on the updated node, edge updates are calculated using attention mechanism. (c) *Node Update*: Nodes are updated using edge aggregation. (d) *Global update*: Updated node attributes are passed through an MLP to predict the conditional intensity, type of next event, and time of next event. (e) *Propagation*: Current node attributes are again fed back to the model as past history when a new event is observed. The red and green lines are the observed events sequentially appearing in time.

Edge Aggregation: a'_k and b'_k are then used by the aggregation function $\rho^{e \rightarrow v}$ to aggregate the edges which have node i as the receiver into $\bar{e}'_i \in \mathbb{R}^{d_e}$. The scalar terms a'_k are normalized to obtain the attention scores, which are used as the weights for weighted element-wise summation of b'_k .

$$\bar{e}'_i = \rho^{e \rightarrow v}(E'_i) := \frac{1}{\sum_{k:r_k=i} a'_k} \sum_{k:r_k=i} a'_k b'_k \quad (17)$$

Multi-headed Attention: The edge update and aggregation steps can be performed independently N_h times simultaneously on the same input. We observed this improves the performance of the model similar to the results observed by Veličković et al. [9] as it allows the network to jointly attend to input projections in various representation subspaces [15]. All the N_h different edge updates $\bar{e}'_i{}^{th}$ corresponding to different attention heads are concatenated into one single vector $\tilde{e}'_i \in \mathbb{R}^{N_h \times d_e}$. $\tilde{e}'_i = [\bar{e}'_i{}^1, \dots, \bar{e}'_i{}^{N_h}]$

Node Update: Node attributes $v_i \in \mathbb{R}^{d_v}$ are updated by passing the multi-headed edge aggregation \tilde{e}'_i through an MLP NN_v .

$$v'_i = \phi^v(\tilde{e}'_i, v_i, \mathbf{u}) := f^v(\{\bar{e}'_i{}^{th}\}_{h=1}^{N_h}) = \text{NN}_v(\tilde{e}'_i) \quad (18)$$

3) *Global Update*: Once the node attributes are updated, we concatenate all the node attributes and pass it through an MLP NN_u to obtain a global attribute

$\mathbf{u}' \in \mathbb{R}^{d_u}$ which is used to predict - conditional intensity for all event types $\lambda_y^* \in \mathbb{R}_+^{|\mathcal{Y}|}$, type of next event $\hat{y}_{i+1} \in \mathcal{Y}$ and time of the next event $\hat{t}_{i+1} \in (t_i, \infty)$. This global attribute represents the history \mathcal{H}_t upto time t_i . We add t_i in subscript to convey the temporal dependence of this global attribute.

$$\mathbf{u}_{t_i} = \mathbf{u}' = \phi^u(\bar{e}, \bar{v}, \mathbf{u}) := \text{NN}_u([\mathbf{v}'_1, \dots, \mathbf{v}'_{|\mathcal{Y}|}]) \quad (19)$$

$$\hat{y}_{i+1} = \arg \max_y \text{Softmax}(\text{NN}_y(\mathbf{u}_{t_i})) \quad (20)$$

$$\hat{t}_{i+1} = \text{NN}_t(\mathbf{u}_{t_i}) \quad (21)$$

We define our model in terms of conditional intensity due to its simplicity. As it is not a probability density, the only restriction is that it has to be non-negative. Moreover, it need not sum up to 1 unlike the conditional density. The conditional intensity of a Hawkes process is defined at every point $t \in \mathbb{R}^+$, however, equation 19 only outputs the hidden representations at the timestamps in the observed sequence. To interpolate the conditional intensity λ^* to the time when an event does not occur, we use the expression proposed by Zuo et al. [5].

$$\lambda^*(t) = \sum_{y=1}^{|\mathcal{Y}|} \lambda_y^*(t) \quad (22)$$

$$\lambda_y^*(t) = \text{Softplus}(\alpha_y \frac{t - t_i}{t_i} + \text{NN}_\lambda(\mathbf{u}_{t_i})_y + \beta_y) \quad (23)$$

The softplus function guarantees that the λ_y^* are non-negative. The first term of equation 23 represents the

contribution of the current event time towards the future and α_y is a hyperparameter. The second term is the history term that encodes how the history \mathcal{H}_t of observed events influence the conditional intensity of a certain event type. The last term β_y incorporates the base intensity of the point process in the absence of events.

Updated node attributes $\{\mathbf{v}'_y\}_{y=1}^{|\mathcal{Y}|}$ now become the previous node attributes $\{\mathbf{v}^y_{t_i}\}_{y=1}^{|\mathcal{Y}|}$ when the next event (t_{i+1}, y_{i+1}) is processed by the model. The entire information flow is illustrated in Fig. 1.

B. Input Embedding

Although RNN allows for a natural ordering in the processing of inputs, as the arrival of the events is not uniform, temporal information needs to be explicitly fed to the model. Directly feeding the time t_i causes issues as the input value increases unbounded or the model may not see any sequence with one specific length which would hurt generalization. To overcome this issue, we follow the positional embedding proposed by Vaswani et al. [15]. The trigonometric functions ensure that the input embedding $\mathbf{x}_{t_i} \in \mathbb{R}^{d_x}$ is bounded and deterministic which enable generalization to longer sequences of unseen lengths.

C. Learning Objectives

Log-Likelihood: We use maximum likelihood estimation (MLE) to learn the parameters of the model. For a sequence $\mathcal{S}_i := \{(t_j, y_j) \mid t_j \leq T, j \in \{1, \dots, L_i\}\}$, the log-likelihood of observing the sequence is given by:

$$\ell(\mathcal{S}_i) = \sum_{j=1}^{L_i} \log \lambda_{y_j}^*(t_j) - \int_0^T \lambda^*(t) dt \quad (24)$$

The first term in the log-likelihood expression Eq. 24, is the log-likelihood of an event occurring at times t_j , we would like this term to be as large as possible. Whereas, the second term signifies the log-likelihood of no events occurring in the times t other than t_j . We would like this term to be as small as possible.

As the sequences \mathcal{S}_i in the dataset are assumed to be i.i.d, the loss function to be minimized is the negative of the sum of the log-likelihood over all sequences $\mathcal{L}_\lambda = -\sum_{i=1}^{N_D} \ell(\mathcal{S}_i)$.

The integral $\Lambda = \int_0^T \lambda^*(t) dt$ in Eq. 24 does not have a closed form solution and needs to be numerically approximated. We use the Monte Carlo estimate [27] given in the supplementary material.

Event Prediction: We are also interested in predicting the type of the next event, we additionally impose a cross-entropy loss term that penalizes when the model mispredicts the type of the next event. For an event (t_j, y_j) , let $\mathbf{y}_j \in \mathbb{R}^{|\mathcal{Y}|}$ be the one-hot encoding of the event type y_j . Hence the next event prediction loss is given by: $\mathcal{L}_y = -\sum_{i=1}^{N_D} \sum_{j=2}^{L_i} \mathbf{y}_j^T \log(\text{NN}_y(\mathbf{u}_{t_j}))$.

Time Prediction: Apart from event prediction, we also want the predicted time of the next event to be

close to the ground truth. Time of the next event is a continuous value that needs to be estimated thus we use an $L2$ penalty to reduce the Mean Squared Error (MSE). The time prediction loss is given by: $\mathcal{L}_t = \sum_{i=1}^{N_D} \sum_{j=2}^{L_i} |t_j - \hat{t}_j|^2$.

V. EXPERIMENTS

A. Datasets

Retweets(RT) [10]: This dataset contains sequences of tweets where each sample has a sequence of tweets of different users.

StackOverflow(SO) [11]: The StackOverflow dataset contains sequences of awards that users were awarded for answering questions on the StackOverflow website. The markers for the events are the various different awards.

MIMIC-II [12]: MIMIC-II dataset contains the visitation of various patients to a Hospital’s ICU. Each sample represents a patient and the events markers are the diagnosis.

Financial Transactions [3]: This dataset contains the transaction records for multiple stocks on a single day. The different events are buy and sell orders at various timestamps.

StarCraft II(SC-II): We introduce a new dataset for benchmarking various discrete event models. Each sequence is a “build-order” which represents a temporally ordered list of various buildings built by the players over the course of a game. The strategy of the players is the underlying complex stochastic process that generates this data. Each sequence in the dataset is an entire Protoss vs Protoss game. Additional details are presented in the supplementary material.

B. Setup

Here we describe the architectural choices we make for various update function. In a'_k , $\text{NN}_\alpha = \text{NN}_\beta$ is simply a linear projection $W \in \mathbb{R}^{d_v \times d_e}$, $\text{NN}_{\alpha'}$ is a single fully connected (FC) layer with leakyReLU [9] non-linearity. NN_v is a linear projection from $\mathbb{R}^{N_h \times d_e}$ to \mathbb{R}^{d_v} . NN_u uses a single FC layer with ReLU non-linearity whereas NN_y , NN_t and NN_λ are simple linear projections of size $\mathbb{R}^{d_u \times |\mathcal{Y}|}$, $\mathbb{R}^{d_u \times 1}$ and $\mathbb{R}^{d_u \times |\mathcal{Y}|}$ respectively.

We are interested in the three evaluation metrics - Log-likelihood, event class prediction accuracy and event time prediction error. We also evaluate these metrics for - Transformer Hawkes Process (THP) [5], Recurrent Marked Temporal Point Process (RMTTP) [3] and Neural Hawkes Process (NHP)[4] for comparison. For a fair comparison, we optimize these models using the same objective as ours and for each evaluation metric we pick the model parameters which lead to the best performance on the validation set. Additional training details are included in the supplementary material.

C. Likelihood

In this section, we use log-likelihood as an evaluation metric as done in previous works [7], [6], [5], [4], [3]. A

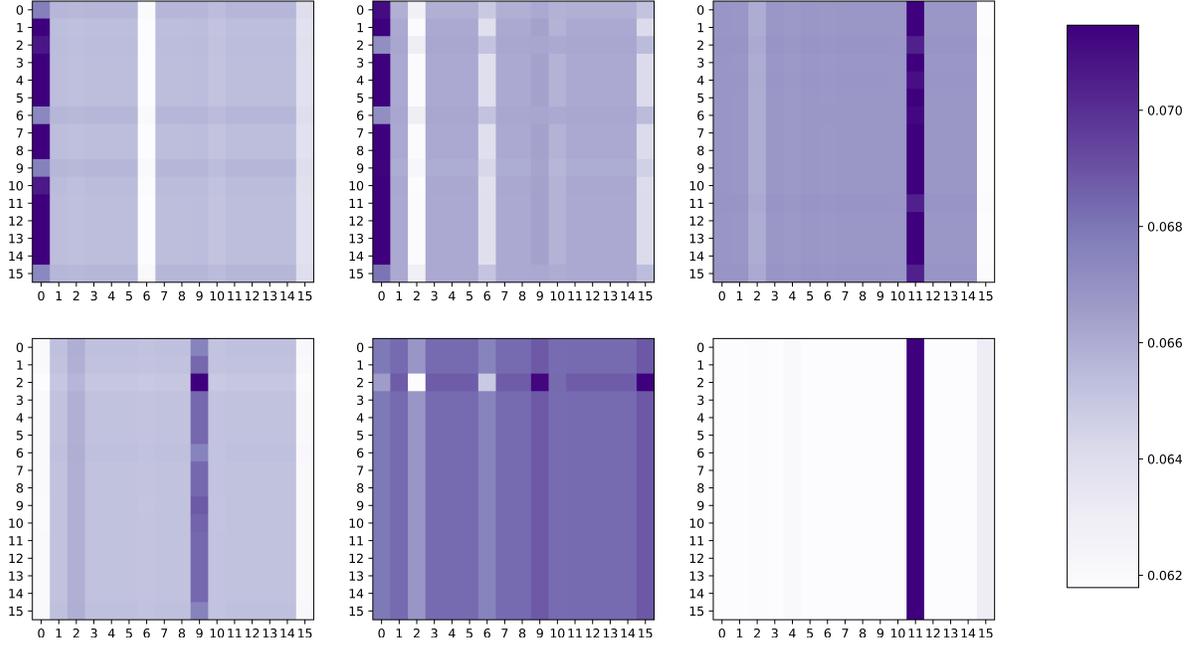


Figure 2: Visualization of various attention heads (row) for different timestamps (column)

Table II: Log-Likelihood Comparison.

Model	RT	SO	MIMIC-II	Financial	SC-II
RMTTPP	-8.8	-0.73	-0.39	-1.71	1.09
NHP	-8.24	-2.41	-1.38	-2.60	-0.81
THP	-7.80	-0.73	0.86	-1.53	1.04
This work	-6.76	-0.56	1.00	-1.26	1.50

Table IV: Time Prediction L_2 Error Comparison.

Model	RT	SO	MIMIC-II	Financial	SC-II
RMTTPP	16899.72	144.34	3.42	26.95	1.48
NHP	17672.54	144.72	3.17	27.88	1.51
THP	16616.13	140.20	0.87	25.75	1.39
This work	15999.68	121.50	1.026	25.37	1.25

Table III: Event Prediction Accuracy Comparison.

Model	RT	SO	Financial	MIMIC-II	SC-II
RMTTPP	49.89	43.49	60.29	59.91	44.25
NHP	48.82	42.81	61.20	76.21	41.14
THP	50.43	43.23	61.67	80.14	43.33
This work	54.42	45.46	61.69	77.34	47.31

higher log-likelihood score would imply that the model can approximate the conditional intensity function $\lambda^*(t)$ well. Table II shows the log-likelihood results. We can see that RGN beats THP in log-likelihood across all the datasets.

D. Event Type Prediction

We are also interested in the predictive performance of the model to predict the type of next event. To improve the performance of the model on this evaluation metric, we added an extra event type prediction loss \mathcal{L}_y in addition to the standard log-likelihood objective. The results are presented in table III. Unlike the log-likelihood results, we see that although RGN has better predictive performance than THP on most of the datasets, it performs worse on the MIMIC-II dataset.

E. Event Time Prediction

To improve the performance of the model on this evaluation metric, we added an extra event time prediction

loss \mathcal{L}_t in addition to the standard log-likelihood objective. Table IV shows the Root Mean-Squared Error (RMSE) of the proposed approach compared to the baselines. The lower the RMSE, the better the performance of the model. We observe that RGN outperforms THP in four of the five datasets however THP again performs better on the MIMIC-II dataset.

F. Goodness of Fit

We would like to verify that the model describes the structure present in the data accurately however it is difficult to work with non-stationary and history dependent distributions present in real world data. The Time-Rescaling Theorem [28] states that any point process with a conditional intensity function can be mapped to a Poisson Process with unit parameter. Alongside the Time-Rescaling theorem, we can use the learnt conditional intensity function to obtain transformed variables z_j which are independent and should be exponentially distributed with rate 1. Since the cumulative distribution function (CDF) of the target distribution is known, we can use P-P plots to measure the deviation from this ideal behaviour. P-P plots plot the empirical CDF with the actual CDF which should be a straight line equally inclined to both the axes. Fig. 3 shows the P-P plots for our Recurrent Graph

Network (RGN) (above) and Transformer Hawkes Process (THP) (below) for three different datasets. It can be seen that for all the three cases, the P-P plot for THP has substantially larger deviations from the expected straight line whereas Recurrent Graph Network (RGN) remains close to the straight line. This shows that RGN better learns the structure in the data compared to state-of-the-art THP.

G. Model Interpretability

The ability to incorporate structure using Recurrent Graph Network helps in making the deep learning model more interpretable. The formulation of our model allows us to see the event-class dependency i.e how a certain class i depends on a certain class j whenever an event occurs allowing us to study the dynamics of inter-class influence as the event stream is observed. Moreover, incorporating multi-headed mechanism allows the model to attend to different information in different subspaces. Fig 2 shows the visualization of two attention heads in the model at three different event timestamps for the StarCraft II dataset. A lighter intensity corresponds to lower attention while a more pronounced color corresponds to higher attention score to a node. Other transformer based models [5], [6] have to average out the attention scores over all events to find attention over event classes thus losing out on crucial temporal information. Although the model allows us to study inter-class dependency at event timestamps, we can easily extend it to continuous time by using spline-based interpolation [29] or other interpolation techniques.

H. Ablation Studies

We conduct ablation experiments to understand the impact of number of attention heads in the multi-headed attention mechanism of Graph Attention Network. We keep the sizes of all the parameters of the model fixed, and only vary the number of attention heads. The evaluation metrics of interest are Log-likelihood, Event type prediction accuracy and Event time prediction error. In addition, we also study the case where the GAT module is completely removed to study the improvement that inclusion of the GAT module brings over a fully LSTM based model.

The result for the experiment on the StarCraft II dataset is shown in table V. As shown in figure 2, various attention heads in the model attend to different features and table V confirms our hypothesis that doing so improves the performance of the model especially on the log-likelihood metric. We also notice a stark drop in model performance when the graph attention module is removed.

I. Reduction in Computational Complexity

One of the major drawbacks of Transformer based approaches is the quadratic order complexity in space and time. For a sequence with length N , the number of activations required to store the self-attention matrix is in the order of $\mathcal{O}(N^2)$. This problem is further exacerbated

when Transformer layers are stacked or multiple heads are used. Another issue that arises during implementation is due to sequence padding of the mini-batch. All the sequences in the mini-batch are padded with zeros to ensure they have the same length as the longest sequence in the mini-batch. All the sequence attention matrices scale as the square of the longest sequence length in the mini-batch $\mathcal{O}(\max(N)^2)$ leading to a lot of wasted memory and redundant computation. In contrast, the attention mechanism in RGN attends over event types rather than individual events. This makes the attention matrix scale in the order of $\mathcal{O}(|\mathcal{Y}|^2)$ instead of $\mathcal{O}(\max(N)^2)$ which leads to dramatic savings in memory consumed as usually $|\mathcal{Y}| \ll \max(N)$. This also makes the memory consumed by model independent of the sequence length. Table VI shows the number of activations in the attention mechanism (in Millions) and total number of operations (MFLOP) for both the aforementioned models. We report the total number of operations performed by both models to find the historical embedding for all events when a sequence with average length is given from all the datasets. The model hyper-parameters are described in the supplementary material. We can see that RGN has lower attention activations and computations compared to THP for all datasets except for MIMIC-II where the trend is reversed. This is due to the fact that the assumption of $|\mathcal{Y}| < N$ fails to hold for small sequences of average length 4 and a large number of event classes 75 present in the dataset.

Table VII shows the amount of GPU memory used by the models. We can see that for the Financial Transactions dataset, only a mini-batch of size 1 can be fit into an Nvidia RTX2080Ti GPU with 11 Gigabytes of memory. No such issues were faced for RGN which has a very small footprint due to the small number of event types - 2 in the dataset. The trend reverses for the MIMIC-II dataset where the large number of classes cause RGN to consume more memory than THP.

VI. CONCLUSION

In this paper, we present a Recurrent Graph Network to learn the underlying Marked Temporal Point Process from event streams. Our key innovation was incorporating structural information using a dynamic attention mechanism over event types rather than all past events in the sequence. This leads to improved performance on log-likelihood and event prediction metrics. Moreover, it also leads to reduction in time and space complexity across almost all datasets. The model improves interpretability by allowing us to understand the influence of one event

Table V: Sensitivity to number of Attention Heads

# Attention Heads	LL	Type Accuracy	Time Error
0	1.03	43.95	1.269
1	1.38	47.00	1.261
2	1.44	47.15	1.259
4	1.49	47.18	1.260
8	1.50	47.31	1.254

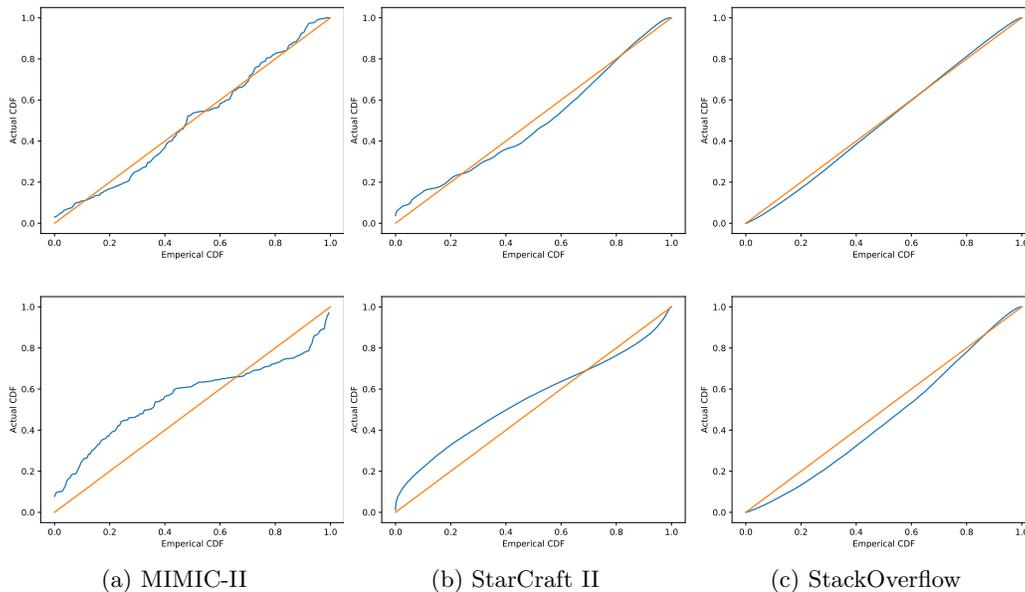


Figure 3: P-P plots for various datasets for this work (above) and THP (below).

Table VI: Complexity Analysis

Dataset	Model	# Computations (MFLOP)	# Activations (M)
RT	THP	89.88	0.41
	This Work	84.10	0.17
SO	THP	3209.67	3.03
	This Work	409.23	1.83
Financial	THP	103294.46	174.58
	This Work	1343.42	32.62
MIMIC II	THP	2.00	0.017
	This Work	23.21	0.081
SC II	THP	1621.84	1.49
	This Work	318.58	1.17

Table VII: Peak GPU Memory Usage (GB).

Dataset	Batch Size	THP	This Work
Retweets	64	2.08	1.03
StackOverflow	16	9.75	3.65
MIMIC-II	64	<u>1.63</u>	5.07
Financial	1	7.15	1.02
StarCraft II	16	10.29	3.63

type over the others as more events are observed by the model. We also present a new interesting benchmarking dataset for point process evaluation.

REFERENCES

- [1] D. Cox and V. Isham, *Point Processes*. Routledge, Dec. 2018.
- [2] A. G. Hawkes, “Spectra of some self-exciting and mutually exciting point processes,” *Biometrika*, 1971.
- [3] N. Du *et al.*, “Recurrent marked temporal point processes: Embedding event history to vector,” pp. 1555–1564, 08 2016.
- [4] H. Mei and J. M. Eisner, “The neural Hawkes process: A neurally self-modulating multivariate point process,” in *NeurIPS 2017*.
- [5] S. Zuo *et al.*, “Transformer Hawkes process,” in *ICML 2020*.
- [6] Q. Zhang *et al.*, “Self-attentive Hawkes process,” in *ICML 2020*.
- [7] O. Shchur *et al.*, “Intensity-free learning of temporal point processes,” in *ICLR 2020*.
- [8] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, 1997.
- [9] P. Veličković *et al.*, “Graph attention networks,” in *ICLR*, 2018.
- [10] Q. Zhao *et al.*, “Seismic,” *SIGKDD 2015*.
- [11] J. Leskovec and A. Krevl, “SNAP Datasets: Stanford large network dataset collection.” <http://snap.stanford.edu/data>, 2014.
- [12] A. E. Johnson *et al.*, “Mimic-iii, a freely accessible critical care database,” *Scientific Data*, 2016.
- [13] Y. Bengio *et al.*, “Learning long-term dependencies with gradient descent is difficult,” *IEEE Transactions on Neural Networks*, 1994.
- [14] R. Pascanu *et al.*, “On the difficulty of training recurrent neural networks,” in *ICML 2013*.
- [15] A. Vaswani *et al.*, “Attention is all you need,” in *NeurIPS 2017*.
- [16] D. Rezende and S. Mohamed, “Variational inference with normalizing flows,” in *ICML 2015*.
- [17] X. Chang *et al.*, *Continuous-Time Dynamic Graph Learning via Neural Interaction Processes*. 2020.
- [18] W. Liang and W. Zhang, “Learning social relations and spatiotemporal trajectories for next check-in inference,” *IEEE TNNLS*, 2020.
- [19] D. Daley and D. Vere-Jones, *An Introduction to the Theory of Point Processes*. Springer-Verlag, 2003.
- [20] J. G. Rasmussen, “Lecture notes: Temporal point processes and the conditional intensity function,” 2018.
- [21] J. Chung *et al.*, “Empirical evaluation of gated recurrent neural networks on sequence modeling,” 2014.
- [22] I. Sutskever *et al.*, “Sequence to sequence learning with neural networks,” in *NeurIPS 2014*.
- [23] K. Cho *et al.*, “Learning phrase representations using RNN encoder-decoder for statistical machine translation,” *CoRR*, 2014.
- [24] N. Srivastava *et al.*, “Unsupervised learning of video representations using lstms,” in *ICML 2015*.
- [25] P. Saha *et al.*, “Physics-incorporated convolutional recurrent neural networks for source identification and forecasting of dynamical systems.”
- [26] P. W. Battaglia *et al.*, “Relational inductive biases, deep learning, and graph networks,” 2018.
- [27] C. P. Robert and G. Casella, *Monte Carlo Statistical Methods (Springer Texts in Statistics)*. Springer-Verlag, 2005.
- [28] F. Papangelou, “The conditional intensity of general point processes and an application to line processes,” *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 1974.
- [29] Hsieh Hou and H. Andrews *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 26, no. 6, pp. 508–517, 1978.