
Conditional Generative Modeling for High-dimensional Marked Temporal Point Processes

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent advancements in generative modeling have made it possible to generate
2 high-quality content from context information, but a key question remains: how to
3 teach models to know when to generate content? To answer this question, this study
4 proposes a novel event generative model that draws its statistical intuition from
5 marked temporal point processes, and offers a clean, flexible, and computationally
6 efficient solution for a wide range of applications involving the generation of asyn-
7 chronous events with high-dimensional marks. We use a conditional generator that
8 takes the history of events as input and generates the high-quality subsequent event
9 that is likely to occur given the prior observations. The proposed framework offers
10 a host of benefits, including considerable representational power to capture intricate
11 dynamics in multi- or even high-dimensional event space, as well as exceptional
12 efficiency in learning the model and generating samples. Our numerical results
13 demonstrate superior performance compared to other state-of-the-art baselines.

1 Introduction

14
15 Generating future events is a challenging yet fascinating task, with numerous practical applications
16 [2, 9, 16, 31]. For instance, a news agency may need to generate news articles in a timely manner, taking
17 into account the latest events and trends. Similarly, an online shopping platform may aim to provide highly
18 personalized recommendations for products, services, or content based on a user’s preferences and behavior
19 patterns over time, as shown in Figure 1. These types of applications are ubiquitous in daily life, and
20 the related data typically consist of a sequence of events that denote when and where each event occurred,
21 along with additional descriptive information such as category, volume, and even text or image, commonly
22 referred to as “marks”. Recent improvements in generative modeling have made it possible
23 to generate high-quality content from contextual information such as language descriptions. However, it
24 remains an open question: how to teach these models to determine the appropriate timing for generating
25 such content based on the history of events.
26
27
28
29
30
31
32
33
34
35

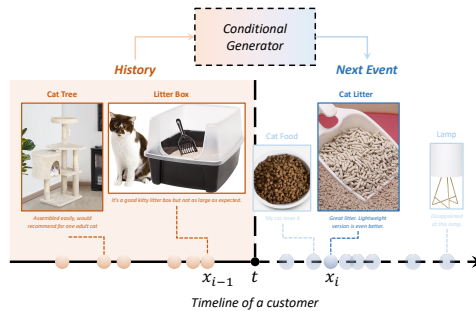


Figure 1: An example of generating high-dimensional content over time. In this example, the conditional generator explores the customer’s next possible activity, including not only the purchase time, but also the item, and even its image or review. The observed events from the customer’s past purchases are represented by yellow dots, while the next generated event is indicated by a blue dot.

36 Point processes have been a popular tool for modeling and generating asynchronous and discrete event
37 data. With the rise of complex systems, advanced neural point processes [6, 18, 25] are proposed as
38 powerful methods to model and simulate data by capturing complex dependencies among observed
39 events. However, due to the use of neural networks, the model likelihood is often analytically
40 intractable, requiring complex and expensive approximations during learning. More seriously, these

41 models face significant limitations in *generating events with high-dimensional marked information*,
 42 as the event simulation relies heavily on the thinning algorithm [20], which can be costly or even
 43 impossible when the mark space is high-dimensional. This significantly restricts the applicability of
 44 these models to modern applications [30, 34], where event data often come with high-dimensional
 45 marks, such as texts and images in police crime reports or social media posts.

46 To tackle these challenges, this paper introduces a novel combination of generative framework and
 47 marked temporal point processes for efficient modeling and generation of high-quality asynchronous
 48 events with high-dimensional marks. The effectiveness of our model is rooted in the ability to
 49 approximate the underlying high-dimensional data distribution through generated samples by a
 50 conditional generator, which takes the history of events as its input. The event history is summarized
 51 by a recurrent neural architecture, allowing for flexible selection based on the application’s needs.
 52 The benefits of our model can be summarized by:

- 53 1. Our model is capable of handling time-stamped high-dimensional marks such as images or texts,
 54 leveraging the power of generative models within the framework of marked point processes;
- 55 2. Our model possesses superior representative power, as it does not confine the conditional intensity
 56 or probability density of the events to any specific parametric form;
- 57 3. Our model outperforms existing state-of-the-art baselines in terms of estimation accuracy and
 58 generating high-quality event series;
- 59 4. Our model excels in computational efficiency during both the training phase and the event
 60 generation process. In particular, our method needs only $\mathcal{O}(N_T)$ for generating N_T events, in
 61 contrast to the thinning algorithm’s complexity of $\mathcal{O}(N^d \cdot N_T)$, where $N \gg N_T$ and d represents
 62 the event dimension.

63 It is important to note that our proposed framework is general and model-agnostic, meaning that a
 64 wide spectrum of generative models and learning algorithms can be applied within our framework.
 65 We present two possible learning algorithms in the Appendix A

66 2 Methodology

67 2.1 Background: Marked temporal point processes

68 Marked temporal point processes (MTPPs) [23] consist of a sequence of *discrete events* over time.
 69 Each event is associated with a (possibly multi-dimensional) *mark* that contains detailed information
 70 of the event. Let $T > 0$ be a fixed time-horizon, and $\mathcal{M} \subseteq \mathbb{R}^d$ be the space of marks. We denote the
 71 space of observation as $\mathcal{X} = [0, T) \times \mathcal{M}$ and a data point in the discrete event sequence as

$$x = (t, m), \quad t \in [0, T), \quad m \in \mathcal{M},$$

72 where t is the event time and m represents the mark. Let N_t be the number of events up to time $t < T$
 73 (which is random), and $\mathcal{H}_t := \{x_1, x_2, \dots, x_{N_t}\}$ denote historical events. Let \mathbb{N} be the counting
 74 measure on \mathcal{X} , i.e., for any measurable $S \subseteq \mathcal{X}$, $\mathbb{N}(S) = |\mathcal{H}_t \cap S|$. For any function $\phi : \mathcal{X} \rightarrow \mathbb{R}$, the
 75 integral with respect to the counting measure is defined as $\int_S \phi(x) d\mathbb{N}(x) = \sum_{x_i \in \mathcal{H}_t \cap S} \phi(x_i)$. The
 76 events’ distribution in MTPPs can be characterized via the conditional intensity function λ , which is
 77 defined to be the occurrence rate of events in the marked temporal space \mathcal{X} given the events’ history
 78 $\mathcal{H}_{t(x)}$, i.e.,

$$\lambda(x|\mathcal{H}_{t(x)}) = \mathbb{E}(d\mathbb{N}(x)|\mathcal{H}_{t(x)})/dx, \quad (1)$$

79 where $t(x)$ extracts the occurrence time of the possible event x . Given the conditional intensity
 80 function λ , the corresponding conditional probability density function (PDF) can be written as

$$f(x|\mathcal{H}_{t(x)}) = \lambda(x|\mathcal{H}_{t(x)}) \cdot \exp\left(-\int_{[t_n, t(x)) \times \mathcal{M}} \lambda(u|\mathcal{H}_{t(u)}) du\right). \quad (2)$$

81 where t_n denotes the time of the most recent event before time $t(x)$. The point process models can
 82 be learned using maximum likelihood estimation (MLE). See all the derivations in Appendix B

83 2.2 Conditional event generator

84 The main idea of the proposed framework is to use a *conditional event generator* to produce the i -th
 85 event $x_i = (t_{i-1} + \Delta t_i, m_i)$ given its previous $i - 1$ events. Here, Δt_i and m_i indicate the time
 86 interval between the i -th event and its preceding event and the mark of the i -th event, respectively.
 87 Formally, this is achieved by a generator function:

$$g(z, \mathbf{h}_{i-1}) : \mathbb{R}^{r+p} \rightarrow (0, +\infty) \times \mathcal{M},$$

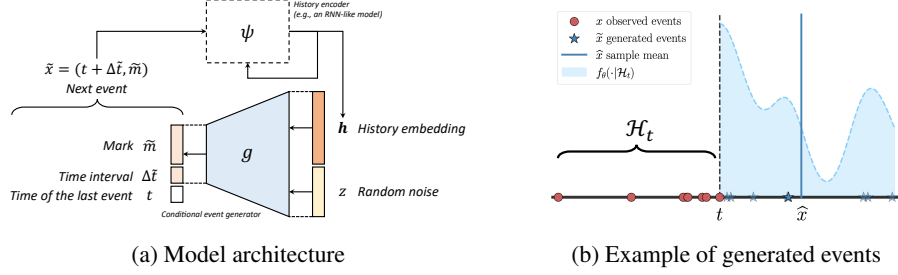


Figure 2: (a) The architecture of the proposed framework, which consists of two key components: A conditional generative model g that generates $(\Delta\tilde{t}, \tilde{m})$ given its history embedding and an RNN-like model ψ that summarizes the events in the history. (b) An example of generated one-dimensional (time only) events $\{\tilde{x}^{(j)}\}$ given the history \mathcal{H}_t . The shaded area suggests the underlying conditional probability density captured by the model with parameters θ .

88 which takes an input in the form of a random noise vector ($z \in \mathbb{R}^r \sim \mathcal{N}(0, I)$) and a hidden
 89 embedding ($\mathbf{h}_{i-1} \in \mathbb{R}^p$) that summarizes the history information up to and excluding the i -th event,
 90 namely, $\mathcal{H}_{t_i} = \{x_1, \dots, x_{i-1}\}$. The output of the generator is the concatenation of the time interval
 91 and mark of the i -th event denoted by $\Delta\tilde{t}_i$ and \tilde{m}_i , respectively. To ensure that the time interval is
 92 positive, we restrict $\Delta\tilde{t}_i$ to be greater than zero.

93 To represent the conditioning variable \mathbf{h}_{i-1} , we use a *history encoder* represented by ψ , which has
 94 a recursive structure such as recurrent neural networks (RNNs) [32] or Transformers [28]. In our
 95 numerical results, we opt for long short-term memory (LSTM) [7], which takes the current event x_i
 96 and the preceding hidden embedding \mathbf{h}_{i-1} as input and generates the new hidden embedding \mathbf{h}_i . This
 97 new hidden embedding represents an updated summary of the past events including x_i . Formally,

$$\mathbf{h}_0 = \mathbf{0} \text{ and } \mathbf{h}_i = \psi(x_i, \mathbf{h}_{i-1}), \quad i = 1, 2, \dots, N_T.$$

98 We denote the parameters of both g and ψ using $\theta \in \Theta$. Figure 2(a) presents the model architecture.

99 **Connection to marked temporal point processes** The proposed framework draws its statistical
 100 inspiration from MTPPs. Unlike other recent attempts at modeling point processes, our framework
 101 *approximates the conditional probability of events using generated samples* rather than directly
 102 specifying the conditional intensity in (1) or PDF in (2) using a parametric model [6, 18, 22, 24, 33].

103 As illustrated by Figure 2(b), when our model generates an event denoted by $\tilde{x} = (t + \Delta\tilde{t}, \tilde{m})$, it
 104 implies that the resulting event \tilde{x} follows a conditional probabilistic distribution that is determined by
 105 the model parameter θ and the event's history \mathcal{H}_t :

$$\tilde{x} \sim f_\theta(x|\mathcal{H}_{t(x)}),$$

106 where f_θ denotes the conditional PDF of the underlying MTPP (2). This design has three main
 107 advantages compared to other point process models:

- 108 1. *Generative efficiency*: The generative nature of our model confers an exceptional efficiency in
 109 simulating a complete event series for any point processes without relying on thinning algorithms
 110 [20]. To exemplify, thinning algorithm (Algorithm 4) has a time complexity of $\mathcal{O}(N^d \cdot N_T)$
 111 to generate N_T events from a history-dependent point process in d -dimensional space \mathcal{X} , with
 112 $N \gg N_T$ being the number of uniformly sampled candidates in one dimension. In contrast, our
 113 generation process (Algorithm 1) only requires a complexity of $\mathcal{O}(N_T)$.
- 114 2. *Expressiveness*: The proposed model enjoys considerable representational power, as it does not
 115 impose any restrictions on the parametric form of the conditional intensity λ or PDF f . The
 116 numerical findings also indicate that our model is capable of capturing complex event interactions,
 117 even in a multi-dimensional space.
- 118 3. *Predictive efficiency*: To predict the next event $\hat{x}_i = (t_{i-1} + \Delta\hat{t}_i, \hat{m}_i)$ given the observed events'
 119 history \mathcal{H}_{t_i} , we can calculate the sample average over a set of generated events $\{\tilde{x}_i^{(l)}\}$ without
 120 the need for an explicit expectation computation, *i.e.*,

$$\hat{x}_i = \int_{(t_{i-1}, +\infty) \times \mathcal{M}} x \cdot f_\theta(x|\mathcal{H}_{t(x)}) dx \approx \frac{1}{L} \sum_{l=1}^L \tilde{x}_i^{(l)},$$

121 where L denotes the number of samples.

Algorithm 1 Event generation process using CEG

Input: Generator g , history encoder ψ , time horizon T
Initialization: $\mathcal{H}_T = \emptyset, \mathbf{h}_0 = \mathbf{0}, t = 0, i = 0$
while $t < T$ **do**
 1. Sample $z \sim \mathcal{N}(0, I)$;
 2. Generate next event $\tilde{x} = (t + \Delta\tilde{t}, \tilde{m})$, where $(\Delta\tilde{t}, \tilde{m}) = g(z, \mathbf{h}_i)$;
 3. $i = i + 1; t = t + \Delta\tilde{t}; x_i = \tilde{x}; \mathcal{H}_T = \mathcal{H}_T \cup \{x_i\}$;
 4. Update hidden embedding $\mathbf{h}_i = \psi(x_i, \mathbf{h}_{i-1})$;
end while
if $t(x_i) \geq T$ **then**
 return $\mathcal{H}_T - \{x_i\}$
else
 return \mathcal{H}_T
end if

122 3 Experiments

123 We evaluate our method using both synthetic and real data and demonstrate the superior performance
124 compared to five state-of-the-art approaches, including (1) Recurrent marked temporal point processes
125 (RMTTP) [6], (2) Neural Hawkes (NH) [18], (3) Fully neural network based model (FullyNN) [22], (4)
126 Epidemic type aftershock sequence (ETAS) [21] model, (5) Deep non-stationary kernel in point process
127 (DNSK) [5]. The first three baselines leverage neural networks to model temporal event data (or only
128 with categorical marks). The last two baselines are chosen for testing multi-dimensional event data.
129 Meanwhile, the DNSK is the state-of-the-art method that uses neural networks for high-dimensional
130 mark modeling. In the following, we refer to our proposed method as the conditional event generator
131 (CEG). Detailed experimental setup and model architectures are presented in Appendix F

132 3.1 Synthetic data

133 We first evaluate our model on synthetic data. To be specific, we generate four one-dimensional
134 (1D) and a three-dimensional (3D) synthetic data sets: Four 1D (time only) data sets include 1,000
135 sequences each, with an average length of 135 events per sequence, and are simulated by two
136 self-exciting processes and two self-correcting processes, respectively, using thinning algorithm
137 (Algorithm 4 in Appendix F). The 3D (time and space) data set also includes 1,000 sequences, each
138 with an average length of 150, generated by a randomly initialized CEG using Algorithm 1

139 To assess the effectiveness of our model in acquiring the underlying data distribution, we computed
140 the mean relative error (MRE) of the estimated conditional intensity and PDF on the testing set, and
141 compared them to the ground truth. Table 1 presents more quantitative results on 1D and 3D data
142 sets, including log-likelihood testing per events and the mean relative error (MRE) of the recovered
143 conditional density and intensity. These results demonstrate the consistent superiority of CEG over
144 other methods across all scenarios. Figure F3 and Figure F4 in Appendix F presents visualizations of
145 the estimated conditional probability density on 1D and 3D synthetic data sets, where CEG accurately
146 captures the complex spatio-temporal point patterns while other baselines fail to do so.

147 3.2 Semi-synthetic data with image marks

148 We test our model’s capability of generating complex high-dimensional marked events on two semi-
149 synthetic data, including time-stamped MNIST (T-MNIST) and CIFAR-100 (T-CIFAR). In these
150 data sets, both the mark (the image category) and the timestamp are generated through a marked
151 point process. Images from MNIST and CIFAR-100 are subsequently chosen at random based
152 on these marks, acting as an high-dimensional representation of the original image category. It’s
153 important to note that during the training phase, categorical marks are excluded, retaining only the
154 high-dimensional images for model learning. Since calculating the log-likelihood for event series with

Table 1: Performance comparison with five baseline methods.

Model	1D self-exciting data			1D self-correcting data			3D synthetic data			3D earthquake data
	Testing ℓ	MRE of f	MRE of λ	Testing ℓ	MRE of f	MRE of λ	Testing ℓ	MRE of f	MRE of λ	Testing ℓ
RMTTP	-1.051 (0.015)	0.437	0.447	-0.975 (0.006)	0.308	0.391	/	/	/	/
NH	-0.776 (0.035)	0.175	0.198	-1.004 (0.010)	0.260	0.363	/	/	/	/
FullyNN	-1.025 (0.003)	0.233	0.330	-0.821 (0.008)	0.322	0.495	/	/	/	/
ETAS	/	/	/	/	/	/	-4.832 (0.002)	0.981	0.902	-3.939 (0.002)
DNSK	-0.649 (0.002)	0.015	0.024	-2.832 (0.004)	0.134	0.185	-2.560 (0.004)	0.339	0.415	-3.606 (0.003)
CEG	-0.645 (0.002)	0.013	0.066	-0.768 (0.005)	0.042	0.075	-2.540 (0.011)	0.049	0.089	-2.629 (0.015)

*Numbers in parentheses present standard error for three independent runs.

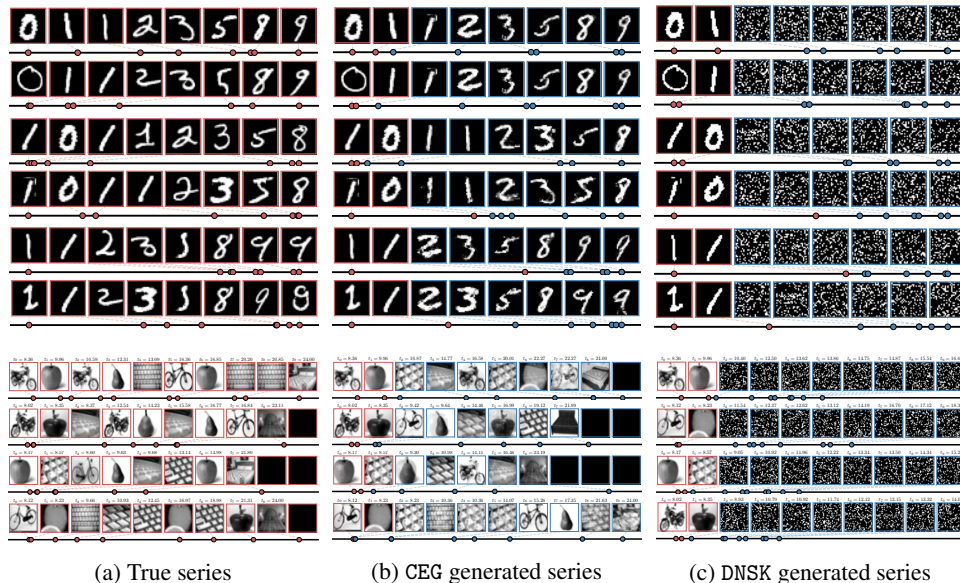


Figure 3: Generated T-MNIST (first row) and T-CIFAR (second row) series using CEG and a neural point process baseline DNSK, with true sequences displayed in the first column. Each event series is generated (blue boxes) given the first two true events (red boxes).

high-dimensional marks is infeasible for CEG (the number of samples needed to estimate density is impractically large), we evaluate the model performance according to: (1) the quality of the generated image marks and (2) the transition dynamics of the entire series. Details of the data generation processes can be found in Appendix F

1. T-MNIST: For each sequence in the data, the actual digit in the succeeding image is the aggregate of the digits in the two preceding marks. The initial two digits are randomly selected from 0 and 1. The digits in the marks must be less than nine. The hand-written image for each mark is then chosen from the corresponding subset of MNIST according to the digit. The time for the entire MNIST series conforms to a Hawkes process with an exponentially decaying kernel.
2. T-CIFAR: The data contains event series that depict a typical day in the life of a graduate student, spanning from 8:00 to 24:00. The marks are sampled from four categories: outdoor exercises, food ingestion, working, and sleeping. Depending on the most recent activity, the subsequent one is determined by a transition probability matrix. Images are selected from the respective categories to symbolize each activity. The activity times follows a self-correcting process.

Figure 3 presents the true T-MNIST series alongside the series generated by CEG and DNSK given the first two events. Our model not only generates high-dimensional event marks that resemble true images, but also successfully captures the underlying data dynamics, such as the clustering patterns of the self-exciting process and the transition pattern of image marks. On the contrary, the DNSK only learns the temporal effects of historical events and struggles to estimate the conditional intensity for the high-dimensional marks. Besides, the grainy images generated by DNSK demonstrate the challenge of simulating credible high-dimensional content using thinning algorithm. This is because the real data points, being sparsely scattered in the high-dimensional mark space, make it challenging for the candidate points to align closely with them in the thinning algorithm.

Similar results are shown in Figure 3 on T-CIFAR data set, where the CEG is able to simulate high-quality daily activities with high-dimensional content at appropriate times. However, the DNSK fails to extract any meaningful patterns from the data, since intensity-based modeling and data generation become ineffectual in high-dimensional mark space.

3.3 Real data

In our real data results, our model demonstrates superior efficacy in generating multi- and high-dimensional event sequences of high quality, which closely resemble real event series.

Northern California earthquake catalog We test our method using the Northern California Earthquake Data [19], which contains detailed information on the timing and location of earthquakes that occurred in central and northern California from 1978 to 2018, totaling 5,984 records with

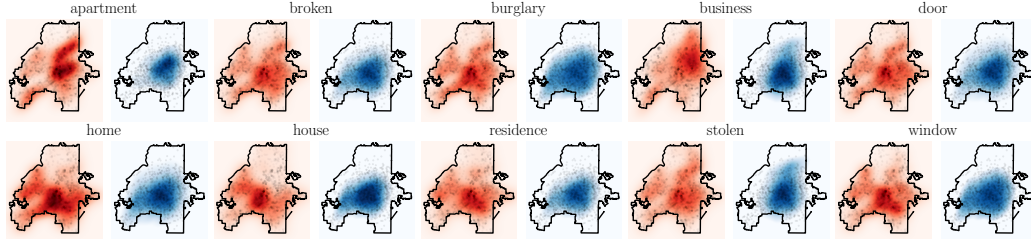
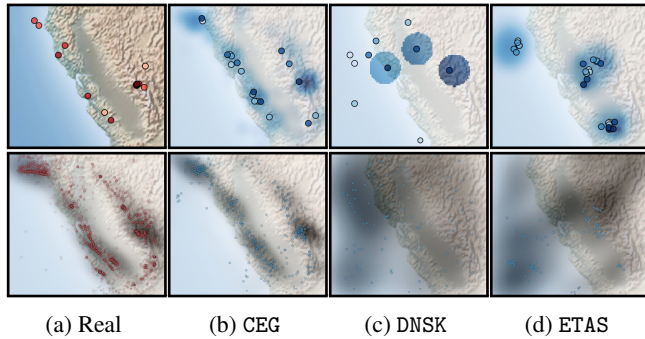


Figure 5: The spatial distributions of the TF-IDF values of 10 crime-related keywords. The heatmap in red and blue represent distributions of TF-IDF value of the keywords in the true and generated events, respectively. The black dots pinpoint the locations of the corresponding events.

188 magnitude greater than 3.5. We divided the data into several sequences by month. In comparison
 189 to other baseline methods that can only handle 1D event data, we primarily evaluated our model
 190 against DNSK and ETAS. we assess the quality of the generated sequences by each model. Our model’s
 191 generation process for new sequences can be efficiently carried out using Algorithm 1, whereas
 192 both DNSK and ETAS requires the use of a thinning algorithm (Algorithm 4) for simulation. We also
 193 compared the estimated conditional probability density functions (PDFs) of real sequences by each
 194 model in Appendix F

195 We compare the generative ability
 196 of each method in Figure 4. The
 197 top left sub-figure features a single
 198 event series selected at random
 199 from the data set, while the rest of
 200 the sub-figures in the first row exhibit
 201 event series generated by each
 202 model, respectively. The quality of
 203 the generated earthquake sequence
 204 using our method is markedly superior
 205 to that generated by DNSK and
 206 ETAS. We also simulate multiple
 207 sequences using each method and
 208 visualize the spatial distribution of
 209 generated earthquakes in the second
 210 row. The shaded area reflects the
 211 spatial density of earthquakes obtained
 212 by KDE and represents the “background
 213 rate” over space. It is evident that
 214 CEG is successful in capturing the
 215 underlying earthquake distribution, while
 216 the two STPP baselines are unable to
 217 do so. Additional results in Figure F6
 218 visualizes the conditional PDF estimated
 219 by CEG, DNSK, and ETAS for an
 220 actual earthquake sequence in testing
 221 set, respectively. The results indicate
 222 that our model is able to capture the
 223 heterogeneous triggering effects among
 224 earthquakes which align with current
 225 understandings of the San Andreas
 226 Fault System [29]. However, both
 227 DNSK and ETAS fail to extract this
 228 geographical feature from the data.



(a) Real (b) CEG (c) DNSK (d) ETAS
 Figure 4: Comparison between real and generated earthquake sequence. The first row displays a single sequence, either real or generated, with the color depth of the dots reflecting the occurrence time of each event. Darker colors represent more recent events. The shaded areas represent the estimated conditional PDFs. The second row shows 1,000 real or generated events, where the gray area indicates the high density of events, which can be interpreted as the “background rate”.

221 **Atlanta crime reports with textual description** We further assess our method using 911-calls-
 222 for-service data in Atlanta. The proprietary data set contains 4644 burglary incidents from 2016 to
 223 2017, detailing the time, location, and a comprehensive textual description of each incident. Each
 224 textual description was transformed into a TF-IDF vector [11], from which the top 10 keywords with
 225 the most significant TF-IDF values were selected. The location combined with the corresponding
 226 10-dimensional TF-IDF vector is regarded as the mark of the incident. We first fit our CEG model
 227 using the preprocessed data, subsequently generate crime event sequences, and then compare them
 228 with the real data.

229 Figure 5 visualizes the spatial distributions of the true and the generated TF-IDF value of each
 230 keyword, respectively, signifying the heterogeneous crime patterns across the city. As we can observe,
 231 our model is capable of capturing such spatial heterogeneity for different keywords and simulating
 232 crime incidents that follow the underlying spatio-temporal-textual dynamics existing in criminological
 233 *modus operandi* [34].

234 **References**

- 235 [1] Akiko Aizawa. An information-theoretic perspective of tf-idf measures. *Information Processing*
236 *& Management*, 39(1):45–65, 2003.
- 237 [2] Homanga Bharadhwaj, Homin Park, and Brian Y Lim. Recgan: recurrent generative adversarial
238 networks for recommendation systems. In *Proceedings of the 12th ACM Conference on*
239 *Recommender Systems*, pages 372–376, 2018.
- 240 [3] Sam Bond-Taylor, Adam Leach, Yang Long, and Chris G Willcocks. Deep generative modelling:
241 A comparative review of vaes, gans, normalizing flows, energy-based and autoregressive models.
242 *IEEE transactions on pattern analysis and machine intelligence*, 2021.
- 243 [4] Xiuyuan Cheng and Hau-Tieng Wu. Convergence of graph laplacian with knn self-tuned kernels.
244 *Information and Inference: A Journal of the IMA*, 11(3):889–957, 2022.
- 245 [5] Zheng Dong, Xiuyuan Cheng, and Yao Xie. Spatio-temporal point processes with deep non-
246 stationary kernels. *arXiv preprint arXiv:2211.11179*, 2022.
- 247 [6] Nan Du, Hanjun Dai, Rakshit Trivedi, Utkarsh Upadhyay, Manuel Gomez-Rodriguez, and
248 Le Song. Recurrent marked temporal point processes: Embedding event history to vector. In
249 *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and*
250 *data mining*, pages 1555–1564, 2016.
- 251 [7] Alex Graves and Alex Graves. Long short-term memory. *Supervised sequence labelling with*
252 *recurrent neural networks*, pages 37–45, 2012.
- 253 [8] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*
254 *in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 255 [9] Ren-Hung Hwang, Yu-Ling Hsueh, and Yu-Ting Chen. An effective taxi recommender system
256 based on a spatio-temporal factor analysis model. *Information Sciences*, 314:28–40, 2015.
- 257 [10] M Chris Jones. Simple boundary correction for kernel density estimation. *Statistics and*
258 *computing*, 3:135–146, 1993.
- 259 [11] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models.
260 *Advances in neural information processing systems*, 34:21696–21707, 2021.
- 261 [12] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
262 *arXiv:1412.6980*, 2014.
- 263 [13] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International*
264 *Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014,*
265 *Conference Track Proceedings*, 2014.
- 266 [14] Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Founda-*
267 *tions and Trends® in Machine Learning*, 12(4):307–392, 2019.
- 268 [15] Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised
269 learning with deep generative models. *Advances in neural information processing systems*, 27,
270 2014.
- 271 [16] Yifei Ma, Balakrishnan Narayanaswamy, Haibin Lin, and Hao Ding. Temporal-contextual rec-
272 ommendation in real-time. In *Proceedings of the 26th ACM SIGKDD international conference*
273 *on knowledge discovery & data mining*, pages 2291–2299, 2020.
- 274 [17] Raghvendra Mall, Rocco Langone, and Johan AK Suykens. Self-tuned kernel spectral clustering
275 for large scale networks. In *2013 IEEE International Conference on Big Data*, pages 385–393.
276 IEEE, 2013.
- 277 [18] Hongyuan Mei and Jason M Eisner. The neural hawkes process: A neurally self-modulating
278 multivariate point process. *Advances in neural information processing systems*, 30, 2017.

- 279 [19] Northern California Earthquake Data Center. UC Berkeley Seismological Laboratory. Dataset.
280 NCEDC, 2014.
- 281 [20] Yosihiko Ogata. On lewis' simulation method for point processes. *IEEE transactions on*
282 *information theory*, 27(1):23–31, 1981.
- 283 [21] Yosihiko Ogata. Space-time point-process models for earthquake occurrences. *Annals of the*
284 *Institute of Statistical Mathematics*, 50:379–402, 1998.
- 285 [22] Takahiro Omi, naonori ueda, and Kazuyuki Aihara. Fully neural network based model for
286 general temporal point processes. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc,
287 E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32.
288 Curran Associates, Inc., 2019.
- 289 [23] Alex Reinhart. A review of self-exciting spatio-temporal point processes and their applications.
290 *Statistical Science*, 33(3):299–318, 2018.
- 291 [24] Oleksandr Shchur, Marin Biloš, and Stephan Günnemann. Intensity-free learning of temporal
292 point processes. In *International Conference on Learning Representations*, 2020.
- 293 [25] Oleksandr Shchur, Ali Caner Türkmen, Tim Januschowski, and Stephan Günnemann. Neural
294 temporal point processes: A review. *arXiv preprint arXiv:2104.03528*, 2021.
- 295 [26] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsu-
296 pervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei,
297 editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of
298 *Proceedings of Machine Learning Research*, pages 2256–2265, Lille, France, 07–09 Jul 2015.
299 PMLR.
- 300 [27] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using
301 deep conditional generative models. *Advances in neural information processing systems*, 28,
302 2015.
- 303 [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
304 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
305 *processing systems*, 30, 2017.
- 306 [29] Robert Earl Wallace. The san andreas fault system, california: An overview of the history,
307 geology, geomorphology, geophysics, and seismology of the most well known plate-tectonic
308 boundary in the world. 1990.
- 309 [30] Alex Williams, Anthony Degleris, Yixin Wang, and Scott Linderman. Point process models for
310 sequence detection in high-dimensional neural spike trains. *Advances in neural information*
311 *processing systems*, 33:14350–14361, 2020.
- 312 [31] David Wilmot and Frank Keller. A temporal variational model for story generation. *arXiv*
313 *preprint arXiv:2109.06807*, 2021.
- 314 [32] Yong Yu, Xiaosheng Si, Changhua Hu, and Jianxun Zhang. A review of recurrent neural
315 networks: Lstm cells and network architectures. *Neural computation*, 31(7):1235–1270, 2019.
- 316 [33] Shixiang Zhu, Haoyun Wang, Xiuyuan Cheng, and Yao Xie. Neural spectral marked point
317 processes. In *International Conference on Learning Representations*, 2022.
- 318 [34] Shixiang Zhu and Yao Xie. Spatiotemporal-textual point processes for crime linkage detection.
319 *The Annals of Applied Statistics*, 16(2):1151–1170, 2022.