

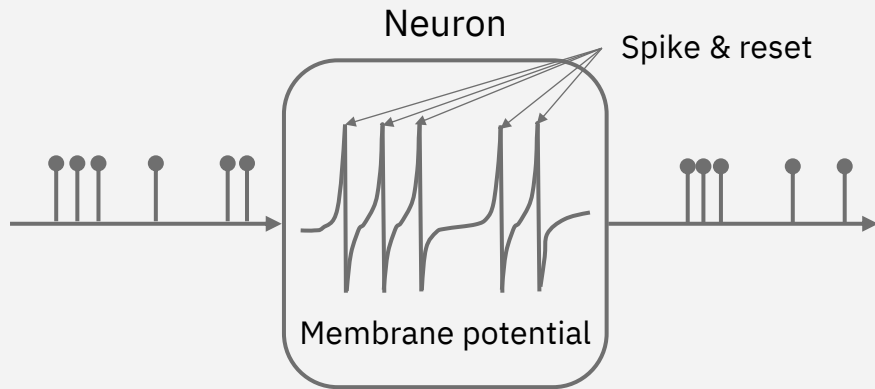
# A Differentiable Point Process with its Application to Spiking Neural Networks

Hiroshi Kajino  
IBM Research - Tokyo

# Opportunities and challenges in spiking neural networks (SNNs)

## Spiking neural networks (SNNs)

- Neurons communicate with spikes
- Neuron has a dynamical state (= membrane potential) evolving in continuous time



## Opportunities

SNNs are often reported to be more energy-efficient

- Spike-based communication
- Efficient information encoding into spikes

## Challenges

- Difficult to train due to discrete nature of spikes
- Time-consuming to simulate SNNs
- Computation time
- Step-size parameter to discretize time axis

# Our approach: Probabilistic model of SNNs

**Point process** (PP) = a probabilistic model of event seq.

$$\mathcal{J}^{\leq t_N} = \{t_1, \dots, t_N\} \subset [0, T]$$

PP is specified by a **conditional intensity function**:

$$\begin{aligned} \lambda(t | \mathcal{J}^{\leq t_n}) dt &= \Pr[t_{n+1} \in [t, t + dt] | \mathcal{J}^{\leq t_n} \text{ and } t_{n+1} \notin (t_n, t)] \\ &= \frac{f(t | \mathcal{J}^{\leq t_n}) dt}{1 - F(t | \mathcal{J}^{\leq t_n})} \end{aligned}$$

↙ Event in  $[t, t + dt]$   
↗ No event in  $(t_n, t)$

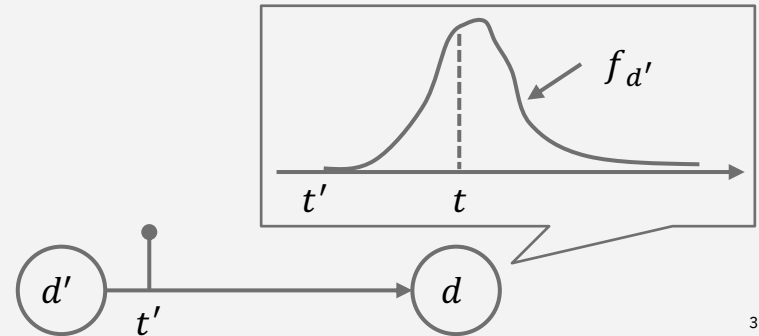
A **probabilistic spiking neuron** is defined as a PP:

- Conditional intensity function  $\propto$  membrane potential

$$\lambda_d(t | \mathcal{J}^{\leq t_n}) = \sigma(u_d(t | \mathcal{J}^{\leq t_n}))$$

- Membrane potential = spike response model

$$u_d(t | \mathcal{J}^{\leq t_n}) = \bar{u}_d + \sum_{(t', d') \in \mathcal{J}^{\leq t_n}} f_{d'}(t - t')$$



# 😊 Probabilistic SNNs can be simulated exactly

**Thinning algorithm** [Lewis+,79] [Ogata,81]

**Input:**  $\lambda(t | \mathcal{T}^{\leq t_n})$ , upperbound  $\bar{\lambda} \geq \lambda(t | \mathcal{T}^{\leq t_n})$

**Output:**  $\mathcal{T}$  Realization of  $\mathcal{PP}(\lambda)$

$k \leftarrow 1, \mathcal{S} \leftarrow \emptyset, \mathcal{T} \leftarrow \emptyset$

Repeat:

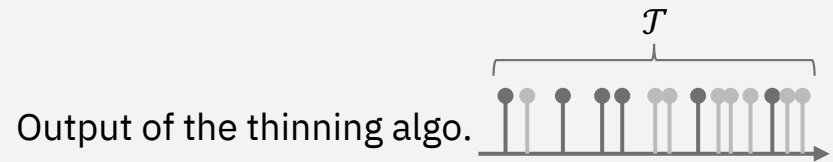
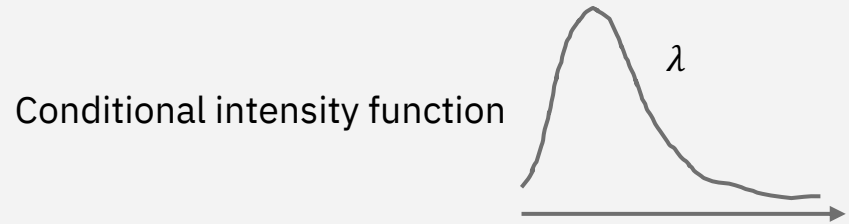
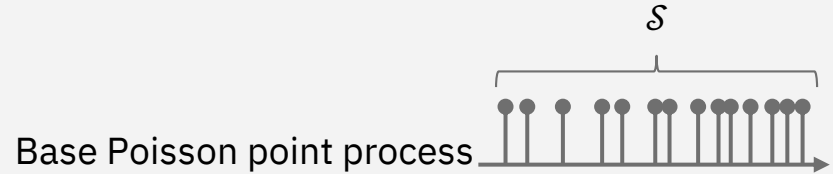
Sample  $s_k \in \mathbb{R} \sim \mathcal{PP}(\bar{\lambda} | \mathcal{S})$

Sample  $d_k \in \{0,1\} \sim \text{Bernoulli}\left(\frac{\lambda(s_k | \mathcal{T})}{\bar{\lambda}}\right)$

If  $d_k = 1$  then  $\mathcal{T} \leftarrow \mathcal{T} \cup \{s_k\}$

$\mathcal{S} \leftarrow \mathcal{S} \cup \{s_k\}$

Return  $\mathcal{T}$



 = Rejected

# 😓 SNNs w/ hidden neurons are difficult to train

**Optimize -ELBO by SGD** [Rezende+,11]  
[Rezende+, 14]

$$\ell(\theta, \phi) = \mathbb{E}_{q(\mathcal{J}_H; \phi)} \left[ \underbrace{-\log p(\mathcal{J}_O, \mathcal{J}_H; \theta) + \log q(\mathcal{J}_H; \phi)}_{\equiv \hat{\ell}(\theta, \phi; \mathcal{J}_O, \mathcal{J}_H)} \right]$$

Observable   Hidden

- $p(\mathcal{J}_O, \mathcal{J}_H; \theta)$ : SNN to be trained
- $q(\mathcal{J}_H; \phi)$ : Variational distribution (point process)

**Existing approach by Rezende+ (2014)**

- Gradient w.r.t.  $\theta$  is easy to estimate
- Gradient w.r.t.  $\phi$  is estimated by

$$\begin{aligned} \frac{\partial \ell}{\partial \phi} &= \frac{\partial}{\partial \phi} \mathbb{E}_{q(\mathcal{J}_H; \phi)} [\hat{\ell}(\theta, \phi; \mathcal{J}_O, \mathcal{J}_H)] \\ &= \mathbb{E}_{q(\mathcal{J}_H; \phi)} [\hat{\ell} \cdot \nabla_{\phi} \log q(\mathcal{J}_H; \phi) + \nabla_{\phi} \hat{\ell}] \end{aligned}$$

whose variance is known to be high ☹

**A remedy: Reparameterization trick**

To make a realization of  $q(\mathcal{J}_H; \phi)$  differentiable



# Our differentiable point process enables the reparameterization trick

## Sampling algorithm of a differentiable PP

**Input:**  $\lambda(t | \mathcal{T}^{\leq t_n})$ , upperbound  $\bar{\lambda} \geq \lambda(t | \mathcal{T}^{\leq t_n})$

**Output:**  $\mathcal{T}$  Realization of  $\partial\mathcal{PP}(\lambda)$

$k \leftarrow 1, \mathcal{S} \leftarrow \emptyset, \mathcal{T} \leftarrow \emptyset$

Repeat:

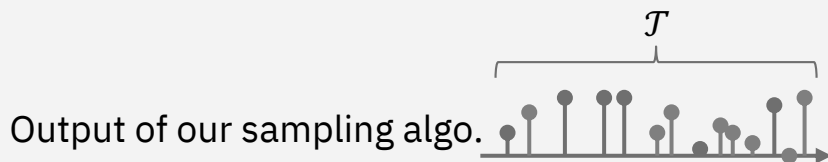
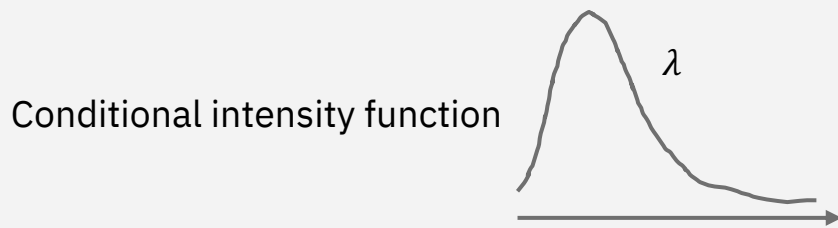
Sample  $s_k \in \mathbb{R} \sim \mathcal{PP}(\bar{\lambda} | \mathcal{S})$

Sample  $d_k \in [0,1] \sim \text{Gumbel\_softmax}_\tau \left( \frac{\lambda(s_k | \mathcal{T})}{\bar{\lambda}} \right)$

$\mathcal{T} \leftarrow \mathcal{T} \cup \{(s_k, d_k)\}$

$\mathcal{S} \leftarrow \mathcal{S} \cup \{s_k\}$

Return  $\mathcal{T}$

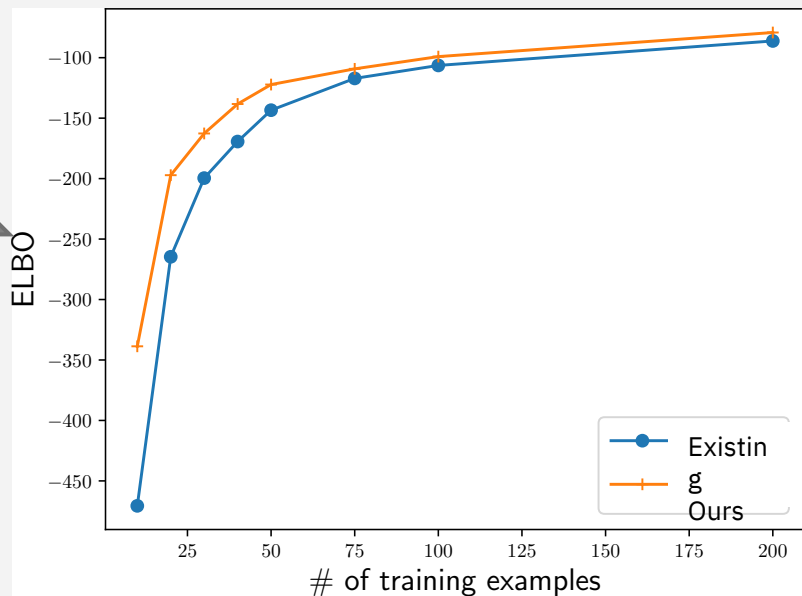


# Empirical studies

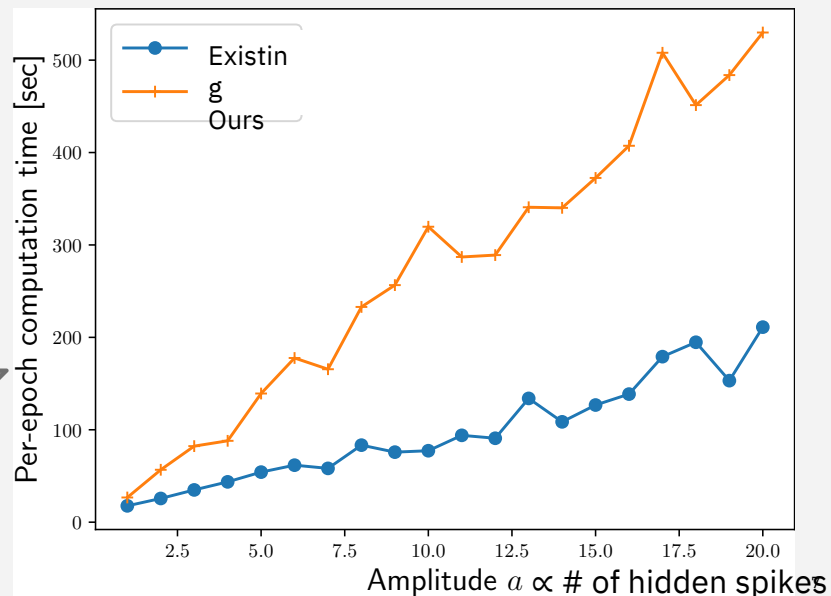
## 1. Standard deviations of gradient estimators

66.3 (Ours) vs. 2,490 (Existing)

## 2. Predictive performance



## 3. Computation time



# Summary

## Our contributions

- A differentiable point process
- A better learning algorithm of probabilistic SNNs

## Findings (vs. the existing work)

- Our gradient estimator has smaller variance
- Smaller var. leads to the better predictive performance
- 2.8x computational overhead

**Code is available under MIT license!**

<https://github.com/ibm-research-tokyo/diffsnn>

