VOWEL: A Local Online Learning Rule for Recurrent Networks of Probabilistic Spiking Winner-Take-All Circuits

Hyeryung Jang

Joint work with Nicolas Skatchkovsky and Osvaldo Simeone

King's College London

ICPR, Jan/12/2021







Machine Learning Today

- Breakthroughs in ML using deep Artificial Neural Networks (ANNs) have come at the expense of massive memory, energy, and time requirements:
 - Energy intensive, particularly for "always-on" edge intelligence applications
 - Many state-of-the-art solutions are not well suited for time-varying and non-stationary environments

Artificial Intelligence / Machine Learni	ng				18	-
Training a sing	le Al model				18 15 15 9 6 6 6 8 15 12 12 12 12 13 14 14 14 14 14 14 14 14 14 14	
as five cars in t	their				Network	
						AlexN
Common carbon footpr	int benchmarks					
in lbs of CO2 equivalent			100000			
Roundtrip flight b/w NY and SF (1 passenger)	1,984	(m)	10000			
Human life (avg. 1 year)	11,023	owe	1000			
American life (avg. 1 year)	36,156	6	100	-		
US car including fuel (avg. 1 lifetime)	126,000		10			
Transformer (213M parameters) w/ neural architecture search	626,155		1	Human	Wa	tsom





Chart: MIT Technology Review - Source: Strubell et al. - Created with Datawrappe

Neuromorphic Computing: Spiking Neural Networks (SNNs)

- Spiking Neural Networks (SNNs) take inspiration from the dynamic, temporally sparse, event-driven learning and inference operations of the human brain
- ANNs to SNNs: From static scalar activation functions to dynamic neuronal models
 - Neurons in the brain sense, process, and communicate over time using sparse binary signals (=spikes)
- Probabilistic SNN Models based on Generalized Linear Model (GLM)
 - More flexible and enable principled differentiable learning rules



Multi-Valued Spikes: Beyond Binary Spikes

- Time-encoded data may assign spikes discrete, categorical, values
 - e.g., polarity (+, -) for data from neuromorphic cameras
- With standard SNN models, this information is either discarded or encoded by increasing the number of signals
 - e.g., producing a binary signal (spike or not), treating as two separate binary signals
- How to model and train an SNN that can process with multi-valued spikes?



Probabilistic (GLM) Winner-Take-All SNNs (WTA-SNNs)

- Discrete data can be encoded via one-hot vector
- It can be processed and produced by Winner-Take-All (WTA) spiking circuit:
 - Group of correlated spiking units, with at most one of the spiking units emitting a spike at any given time



- WTA circuit: categorical one-hot GLM
- Synaptic weights defined by matrices
- Membrane potential

$$\boldsymbol{u}_{i,t} = \sum_{j \in \mathcal{P}_i} \mathbf{W}_{j,i} \vec{\boldsymbol{s}}_{j,t-1} + \mathbf{W}_i \overleftarrow{\boldsymbol{s}}_{i,t-1} + \boldsymbol{\vartheta}_i$$

- Probability of WTA *i* to emit a spike at unit *c* $p_{i,t}(c; \boldsymbol{\theta}_i) = \boldsymbol{\sigma}_{SM}^c(\boldsymbol{u}_{i,t})$
- Log-prob: categorical (neg) cross-entropy $\log p_{i,t}(\boldsymbol{s}_{i,t}; \boldsymbol{\theta}_i) = \overline{H}(\boldsymbol{s}_{i,t}, \boldsymbol{\sigma}_{\text{SM}}(\boldsymbol{u}_{i,t}))$

Training Probabilistic WTA-SNNs

- The visible circuits encode the desired behavior $x_{\leq T}$ of the network as specified by training data
- The spiking behavior $h_{\leq T}$ of the hidden circuits is not specified by data



- Maximum likelihood: $\max_{\theta} \log p_{\theta}$ (visible)
 - Maximize the probability that the model produces the desired outputs for the visible circuits
 - $\log p_{\theta}(\text{visible}) = \log \sum_{\text{hidden}} p_{\theta}(\text{visible, hidden}) \ge \mathbb{E}_{q(\text{hidden})} \left[\log \frac{p_{\theta}(\text{visible, hidden})}{q(\text{hidden})} \right]$
 - Forward distribution log q(hidden) = $\sum_{t} \sum_{i:\text{hidden}} \overline{H}(\boldsymbol{h}_{i,t}, \boldsymbol{\sigma}_{\text{SM}}(\boldsymbol{u}_{i,t}))$
 - We specifically optimize the **lower bound (ELBO)**: $L(\theta) = \sum_{t} \sum_{i:visible} \overline{H}(x_{i,t}, \sigma_{SM}(u_{i,t})) \alpha \cdot reg$

- Use an unbiased estimate (via REINFORCE) of the gradient of the bound $\nabla_{\theta} L(\theta)$
- For visible circuits, we have two-factor rule:

$$\mathbf{W}_{j,i} \leftarrow \mathbf{W}_{j,i} + \eta \cdot \left(\mathbf{x}_{i,t} - \boldsymbol{\sigma}_{\mathrm{SM}}(\boldsymbol{u}_{i,t}) \right) \cdot \left(\vec{\boldsymbol{s}}_{j,t-1} \right)^{\mathsf{T}}$$

post-synaptic pre-synaptic trace predictive error



• For hidden circuits, we have three-factor rule:

$$\mathbf{W}_{j,i} \leftarrow \mathbf{W}_{j,i} + \eta \cdot \sum_{i:\text{visible}} \overline{H}\left(\mathbf{x}_{i,t}, \boldsymbol{\sigma}_{\text{SM}}(\boldsymbol{u}_{i,t})\right) \cdot \left(\mathbf{h}_{i,t} - \boldsymbol{\sigma}_{\text{SM}}(\boldsymbol{u}_{i,t})\right) \cdot \left(\mathbf{\vec{s}}_{j,t-1}\right)^{\mathsf{T}}$$
post-synaptic pre-synaptic trace predictive error

Experiments

- Neuromorphic dataset obtained by filming moving MNIST digits displayed on a screen or moving gestures with a neuromorphic camera
 - MNIST-DVS, DVS-Gesture



- WTA-SNN
 - All WTA circuits consist of two spiking units (for +, signs)
 - Consider a generic, non-optimized, network architecture with H fully connected hidden circuits
- Benchmarks: (binary) probabilistic GLM model and deterministic LIF model
 - With same number of spiking units (= 2H) and per-sign inputs
 - With $\frac{1}{2}$ number of spiking units (= H) with unsigned (discarding the sign) inputs

Results

• MNIST-DVS

• H = 16 hidden circuits (2*H* hidden spiking units)



- WTA-SNN (C = 2) trained with VOWEL outperforms conventional binary SNN solutions (C = 1)
- The capability of WTA circuits to process information encoded both in the spikes' timings and their signs

- DVS-Gesture
 - Test accuracy of VOWEL and DECOLLE (deterministic SNN)

	D · · ·		•	
Model	Period	Н	Accuracy	
	1 ms	2048	$61.42 \pm 2.92 \ \%$	
	1 ms	1024	57.75 ± 3.22 %	
DECOLLE	1 ms	512	56.42 ± 2.03 %	
	10 ms	512	$34.72\pm0.75~\%$	
	20 ms	512	$26.38 \pm 0.28 \ \%$	
VOWEL	20 ms	256	$60.26 \pm 0.91\%$	
	20 ms	128	$57.96\pm0.11\%$	
	20 ms	64	$57.89\pm0.35\%$	

- VOWEL can operate with a smaller number of hidden units and coarser sampling rate
- The ability to directly distinguish patterns encoded in the values of the spikes

Thank you!

Any Questions?

hyeryung.jang@kcl.ac.uk