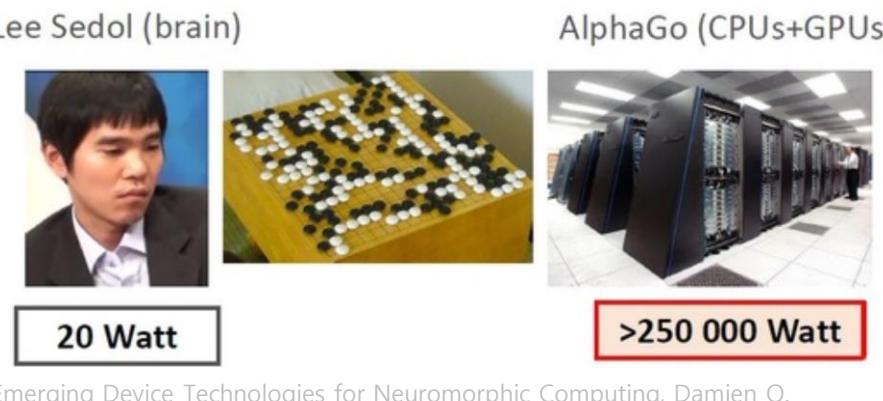


Introduction

Training a single AI model can emit as much carbon as five cars in their lifetimes



Deep learning has a terrible carbon footprint.
Common carbon footprint benchmarks

in lbs of CO₂ equivalent

Roundtrip flight b/w NY and SF (1 passenger)	1,984
Human life (avg. 1 year)	11,023
American life (avg. 1 year)	36,156
US car including fuel (avg. 1 lifetime)	126,000
Transformer (213M parameters) w/ neural architecture search	626,155

<https://www.technologyreview.com/s/613630/training-a-single-ai-model-can-emit-as-much-carbon-as-five-cars-in-their-lifetimes/>

It is important to improve energy efficiency of neural networks!

Publications

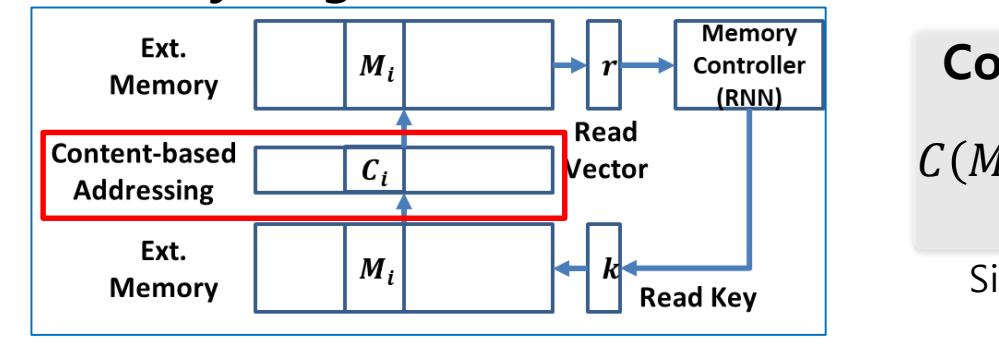
- [J3] Seijo Kim, **Seongsik Park**, Byunggook Na, Jongwan Kim, Sungroh Yoon, "Towards fast and accurate object detection in bio-inspired spiking neural networks through Bayesian optimization," *IEEE Access*, vol. 9, pp. 2633–2643, Jan 2021.
- [J2] **Seongsik Park**, Jaehee Jang, Seijo Kim, Byunggook Na, Sungroh Yoon, "Memory-Augmented Neural Networks on FPGA for Real-Time and Energy-Efficient Question Answering," *IEEE TVLSI*, Vol. 20, no. 1, pp. 162–175, Jan 2021.
- [C8] **Seongsik Park**, Seijo Kim, Byunggook Na, Sungroh Yoon, "T2FSNN: Deep Spiking Neural Networks with Time-to-first-spike Coding," in *Proceedings of Design Automation Conference (DAC)*, 2020.
- [C7] **Seongsik Park**, Jongwan Kim, Sungroh Yoon, "Energy-aware Placement for SRAM-NVM Hybrid FPGAs," in *Proceedings of Design, Automation and Test in Europe (DATE)*, 2020.
- [C6] Seijo Kim, **Seongsik Park**, Byunggook Na, Sungroh Yoon, "Spiking-YOLO: Spiking Neural Network for Energy-efficient Object Detection," in *Proceedings of AAAI Conference on Artificial Intelligence (AAAI)*, 2020.
- [C5] **Seongsik Park**, Seijo Kim, Hyeokjun Choe, Sungroh Yoon, "Fast and Efficient Information Transmission with Burst Spikes in Deep Spiking Neural Networks," in *Proceedings of Design Automation Conference (DAC)*, 2019.
- [C4] **Seongsik Park**, Jahee Jang, Seijo Kim, Sungroh Yoon, "Energy-Efficient Inference Accelerator for Memory-Augmented Neural Networks on FPGA," in *Proceedings of Design, Automation and Test in Europe (DATE)*, 2019.
- [J1] Seil Lee, Hanjoo Kim, **Seongsik Park**, Seijo Kim, Hyeokjun Choe, Sungroh Yoon, "CloudSocket: Fine-Grained Power Sensing System for Datacenters," *IEEE Access*, vol. 6, no. 1, pp. 49601–49610, 2018.
- [C3] **Seongsik Park**, Seijo Kim, Seil Lee, Ho Bae, Sungroh Yoon, "Quantized Memory-Augmented Neural Networks," in *Proceedings of AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [C2] Hyeokjun Choe, Seil Lee, Hyunha Nam, **Seongsik Park**, Seijo Kim, Eui-Young Chung, Sungroh Yoon, "Near-Data Processing for Differentiable Machine Learning Models," in *Proceedings of International Conference on Massive Storage Systems and Technology (MSST)*, 2017.
- [C1] Seil Lee, Hanjoo Kim, **Seongsik Park**, Seijo Kim, Hyeokjun Choe, Chang-Sung Jeong, Sungroh Yoon, "CloudSocket: Smart Grid Platform for Datacenters," in *Proceedings of IEEE International Conference on Computer Design (ICCD)*, 2016.

Topic 1

Energy-efficient DNNs

Memory-augmented Neural Networks + Fixed-point Quantization (AAAI-18)

Memory-augmented Neural Networks



$$\text{Similarity measure function}$$

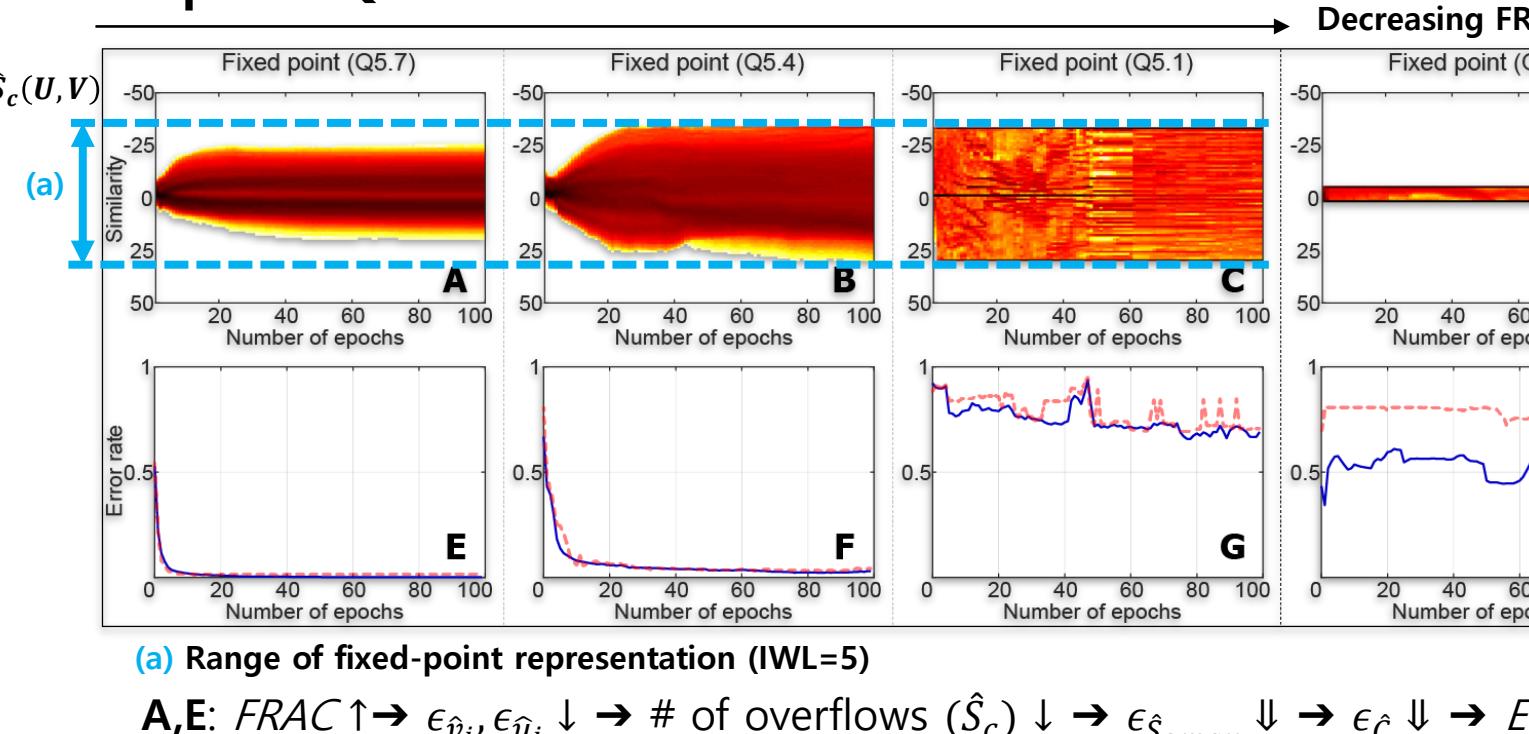
Fixed-point Quantization Error precision

$$|\epsilon_{\hat{u}}| < \begin{cases} 2^{-FRAC} & (\text{if } |u| < 2^{IWL}) \\ 2^{IWL} - |u| & (\text{overflow}) \end{cases}$$

$$\text{overflow}$$

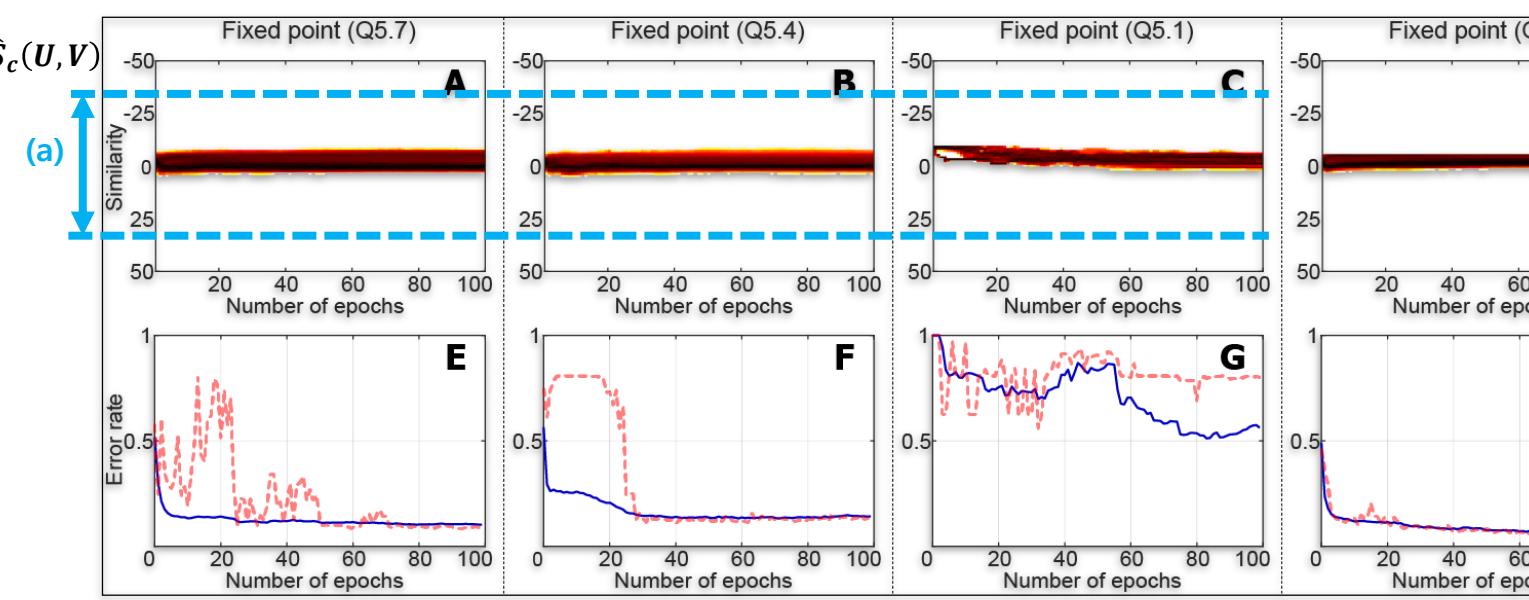
$$\begin{aligned} \text{Fixed point (Q5,2)} & u = 3.875 & \text{Sign} & 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \\ & \hat{u} = 3.75 & \text{Fractional (2bit)} & 1 \ 1 \\ & |\epsilon_{\hat{u}}| = 0.125 & \text{Integer (5bit)} & 0 \ 0 \ 0 \ 0 \ 1 \\ \text{Fixed point (Q1,6)} & u = 3.875 & \text{Sign} & 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \\ & \hat{u} = 1.984375 & \text{Fractional (6bit)} & 1 \ 1 \ 1 \ 1 \ 1 \ 1 \\ & |\epsilon_{\hat{u}}| = 1.890625 & \text{Integer (1bit)} & 1 \end{aligned}$$

Fixed-point Quantization on Conventional MANNs



Proposed Methods - Quantized MANN (Q-MANN)

	Similarity	Similarity measure function
Previous	Cosine similarity (dot product)	$S_c(U, V) \approx \sum_i u_i v_i$
Proposed (Q-MANN)	Hamming similarity (XNOR)	$S_h(\hat{U}, \hat{V}) = \sum_i S_{\hat{u}_i} S_{\hat{v}_i} \sum_{k=0}^{n-2} W_k XNOR(\hat{u}_{ik}, \hat{v}_{ik})$
	Manhattan distance (subtraction)	$\frac{\partial S_h(\hat{U}, \hat{V})}{\partial \hat{u}_i} \approx S_{\hat{u}_i} 2^a (S_{\hat{u}_i} - S_{\hat{v}_i}) - \sum_{k=0}^{n-2} S_{\hat{v}_i} 2^a (\hat{u}_{ik} - \hat{v}_{ik})$

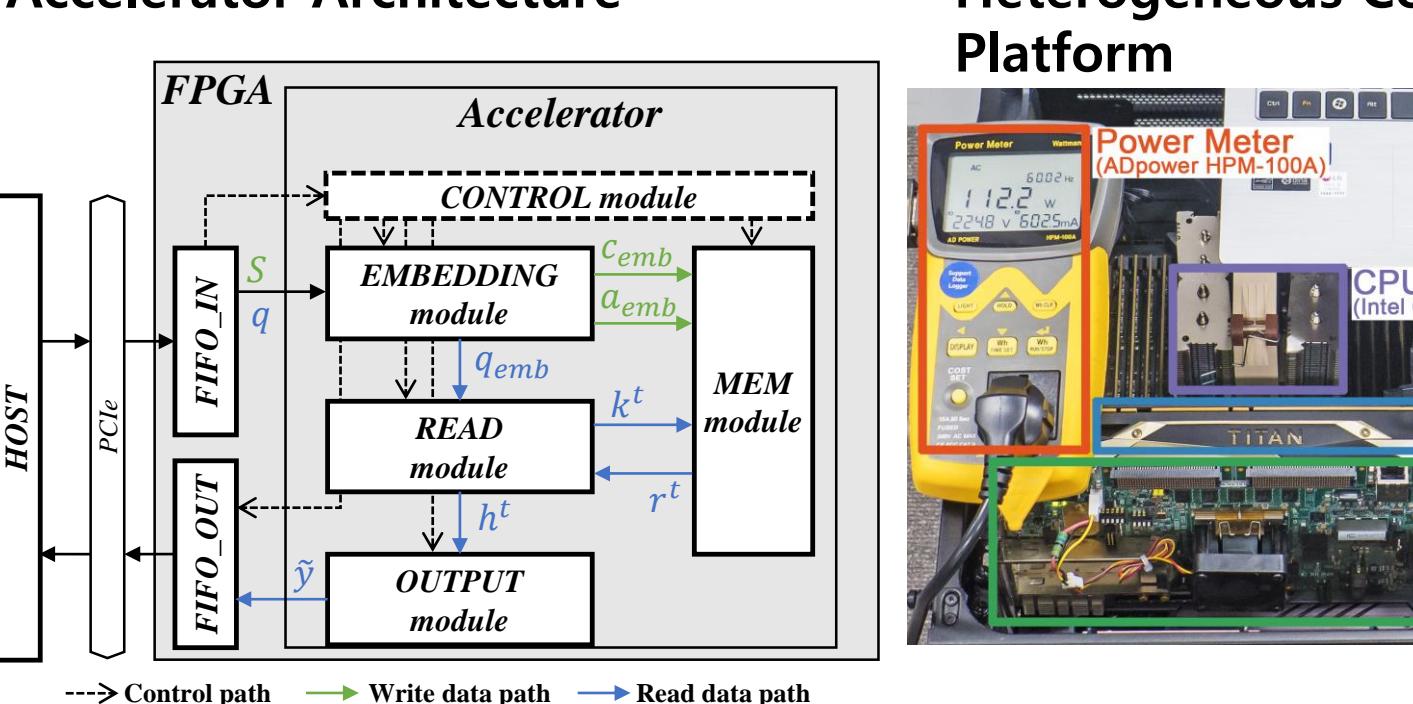


Topic 2

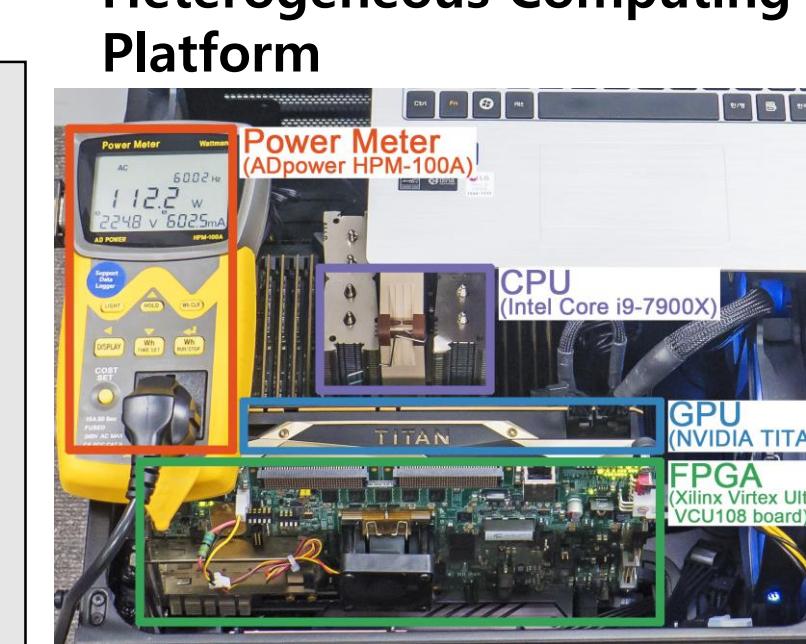
Application-specific inference accelerator

Memory-augmented Neural Networks + Question Answering Accelerator on FPGA (DATE-19, TVLSI-21)

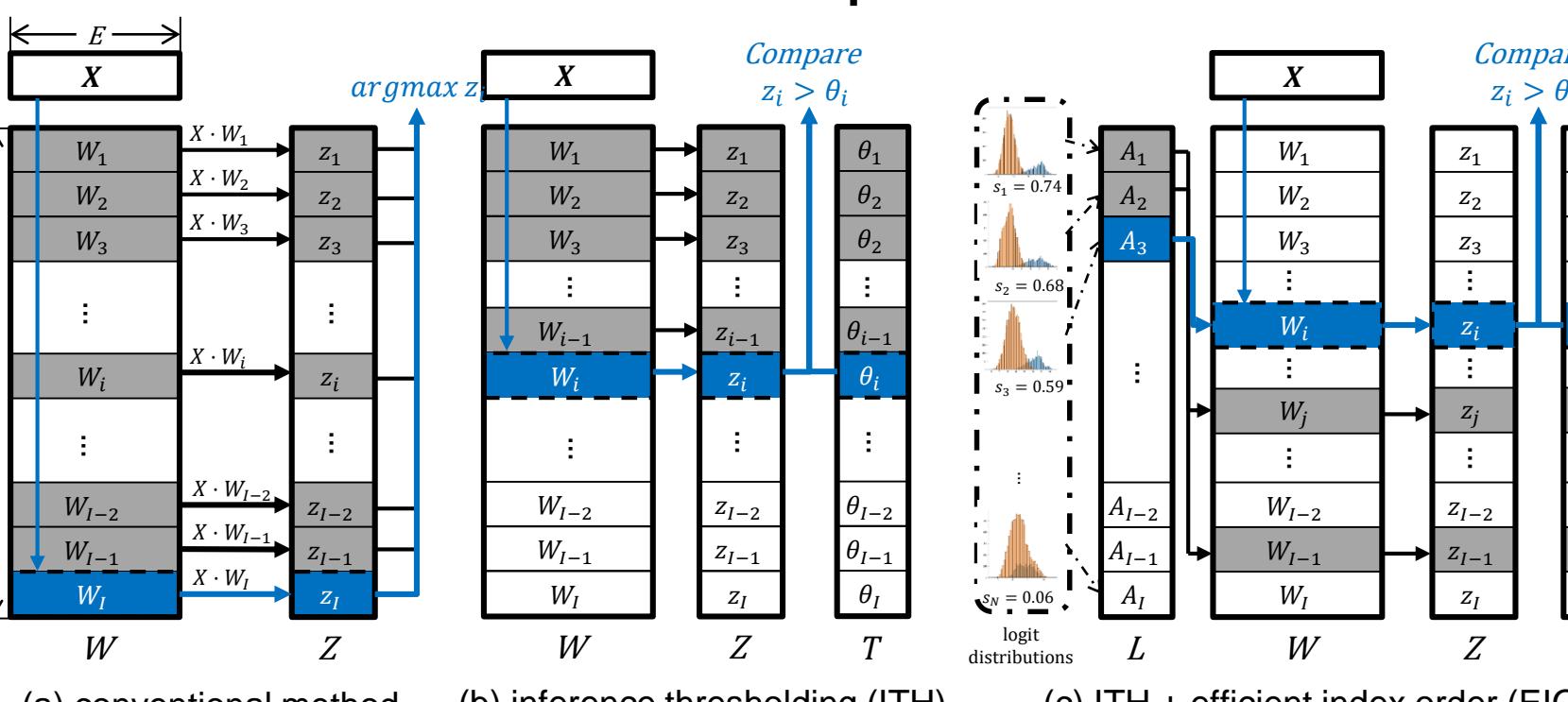
Accelerator Architecture



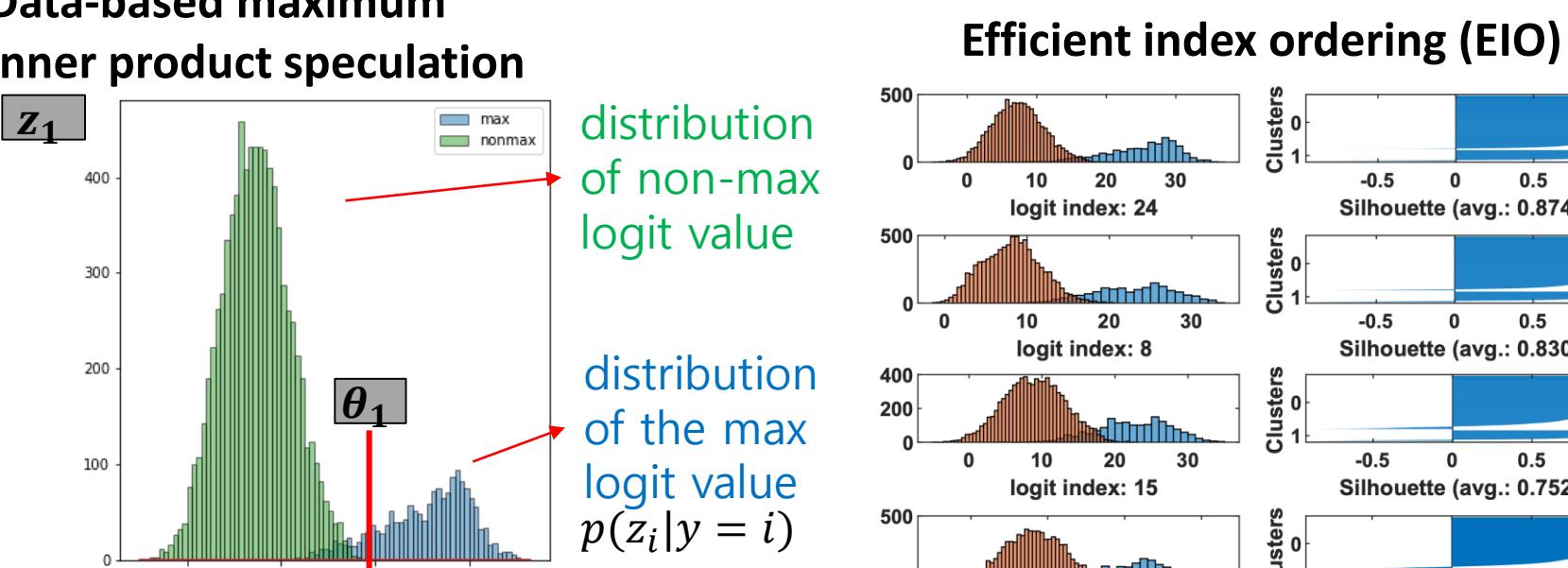
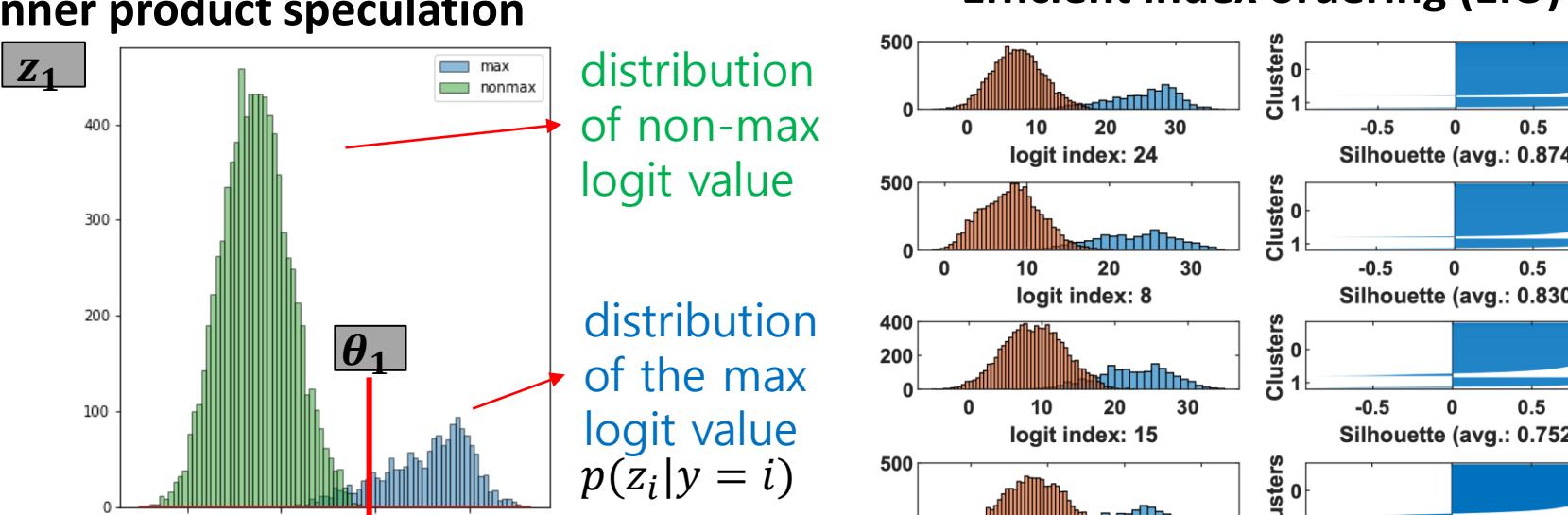
Heterogeneous Computing Platform



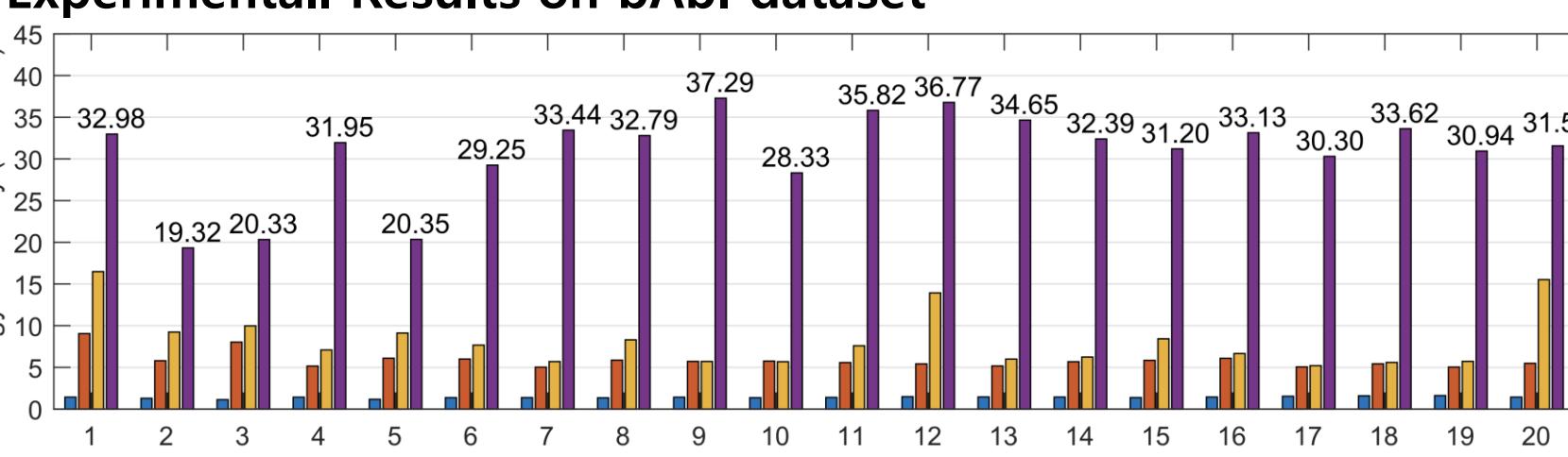
Fast Inference Methods on the Proposed Accelerator



Data-based maximum inner product speculation



Experimental Results on bAbI dataset

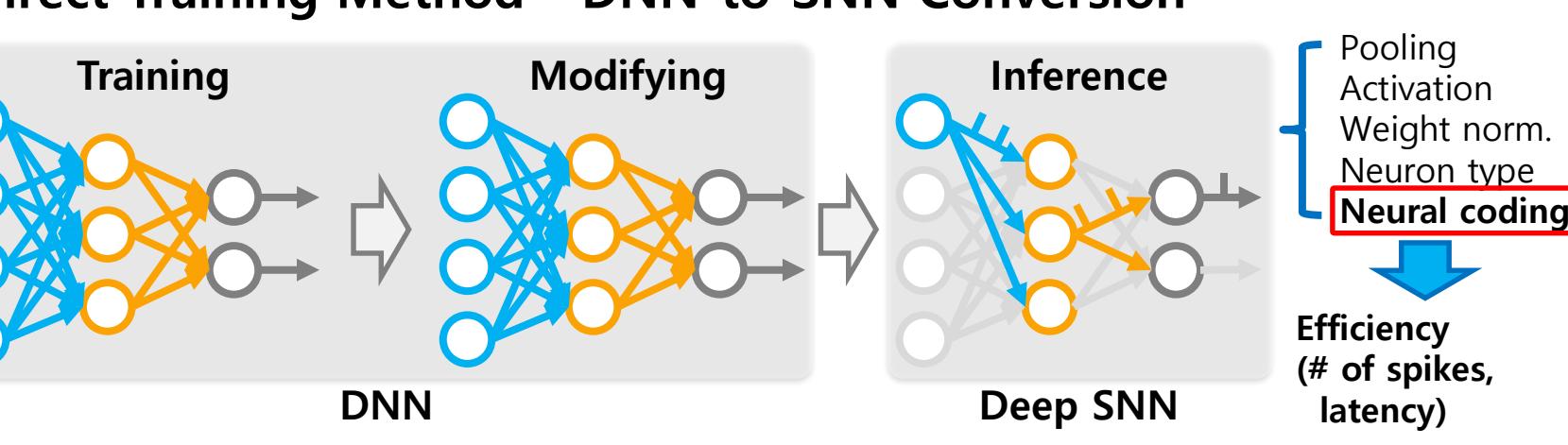


Topic 3

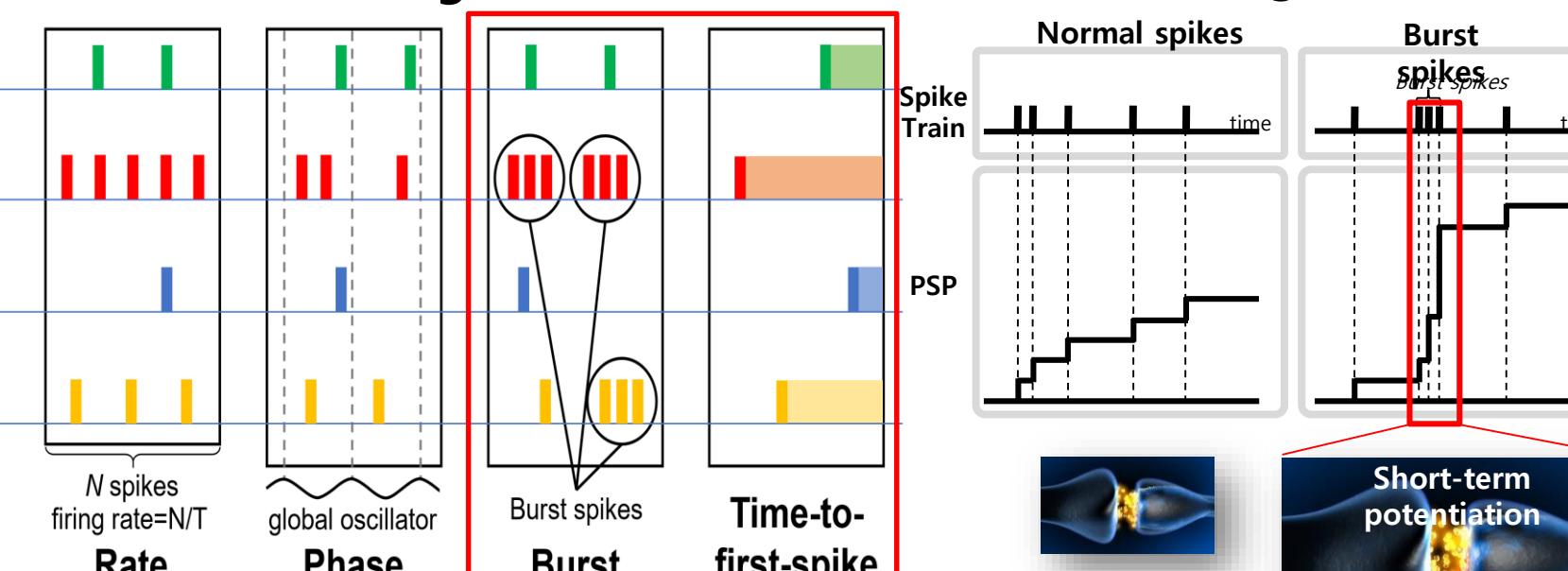
Spiking Neural Networks

Improving Efficiency of Inference on Deep SNNs - Burst coding and Time-to-first-spike coding (DAC-19, DAC-20)

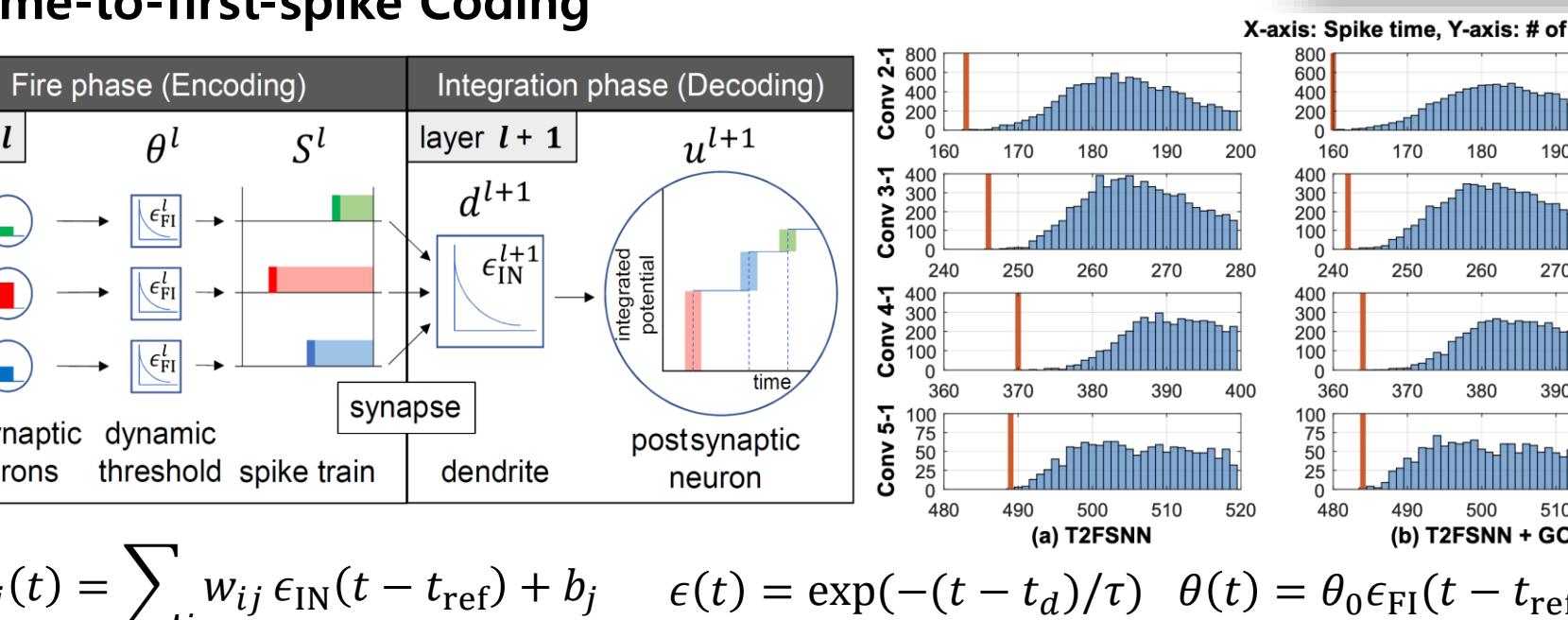
Indirect Training Method - DNN-to-SNN Conversion



Various Neural Coding Schemes



Time-to-first-spike Coding



$$z_j(t) = \sum_i w_{ij} \epsilon_{IN}(t - t_{ref}) + b_j \quad \epsilon(t) = \exp(-(t - t_d)/\tau) \quad \theta(t) = \theta_0 \epsilon_{FI}(t - t_{ref})$$

Experimental Results

