



Neuromorphic Computing: Spiking Neural Networks For Edge AI

Mini T V

Associate Professor, Department of Computer Science, Sacred Heart College (Autonomous), Chalakudy, India.

Article information

Received: 12th December 2025

Received in revised form: 13th January 2026

Accepted: 14th February 2026

Available online: 12th March 2026

Volume: 1

Issue: 1

DOI: <https://doi.org/10.5281/zenodo.18977649>

Abstract

Neuromorphic computing represents a paradigm shift in artificial intelligence by mimicking biological neural networks through spiking neural networks (SNNs). This paper explores neuromorphic computing architectures for energy-efficient AI inference at the edge. We analyze key neuromorphic hardware platforms including Intel Loihi 2, IBM TrueNorth, and BrainScaleS, examining their architectural innovations, spike-timing-dependent plasticity (STDP) learning mechanisms, and event-driven computation models. Performance evaluations demonstrate that SNNs achieve 100-1000× energy efficiency improvements compared to conventional deep neural networks for edge inference tasks. We present implementation strategies for deploying SNNs on resource-constrained edge devices, addressing challenges in spike encoding, temporal dynamics, and neuromorphic algorithm design. Our analysis reveals that neuromorphic computing offers a compelling solution for ultra-low-power AI applications in IoT, robotics, and embedded systems.

Keywords:- Neuromorphic Computing, Spiking Neural Networks, Edge AI, STDP, Event-Driven Computing, Low-Power Inference

I. INTRODUCTION

The exponential growth of edge computing devices and Internet-of-Things (IoT) applications has created an urgent demand for energy-efficient artificial intelligence inference. Traditional deep learning approaches, while highly accurate, consume significant power and require substantial computational resources, making them unsuitable for battery-powered edge devices [1]. Neuromorphic computing emerges as a bio-inspired paradigm that addresses these limitations through event-driven, asynchronous computation modeled after biological neural systems. Spiking Neural Networks (SNNs) form the computational foundation of neuromorphic systems, processing information through discrete spike events rather than continuous activations [2]. This event-driven approach enables remarkable energy efficiency, as computation occurs only when spikes are transmitted between neurons. Recent neuromorphic hardware platforms achieve energy consumption in the range of picojoules per synaptic operation, representing orders of magnitude improvement over conventional GPU-based neural network inference [3].

This paper examines the architecture and implementation of neuromorphic computing systems for edge AI applications. We analyze leading neuromorphic platforms, including Intel Loihi 2 with 1 million neurons and 120 million synapses [4], IBM TrueNorth featuring 4096 neurosynaptic cores [5], and analog platforms like BrainScaleS that operate 10,000× faster than biological real-time [6]. Our investigation encompasses spike encoding mechanisms, temporal learning algorithms, and deployment strategies for real-world edge inference

scenarios. The remainder of this paper is organized as follows: Section II reviews neuromorphic hardware architectures and their design principles. Section III examines spiking neural network models and learning algorithms. Section IV presents performance analysis and energy efficiency metrics. Section V discusses implementation challenges and deployment strategies for edge devices. Section VI concludes with future research directions.

II. NEUROMORPHIC HARDWARE ARCHITECTURES

A. Intel Loihi Architecture

Intel Loihi 2, fabricated in 4nm process technology, represents the state-of-the-art in digital neuromorphic computing [4]. The architecture features 128 neuromorphic cores, each containing 8,192 compartmental neuron models with programmable dynamics. The chip supports asynchronous spike communication through a hierarchical mesh network, enabling scalable inter-core connectivity. Each neuron implements the leaky integrate-and-fire (LIF) model with configurable time constants and refractory periods. The synaptic crossbar architecture in Loihi 2 allows each neuron to connect to up to 64,000 synapses, with 8-bit weight precision and ternary learning rules. On-chip spike-timing-dependent plasticity (STDP) enables autonomous learning without external processor intervention [7]. The chip consumes approximately 300mW during active inference, achieving 4.8 trillion synaptic operations per second with energy efficiency of 0.26pJ per synaptic operation.

B. IBM TrueNorth Architecture

IBM TrueNorth implements a massively parallel neuromorphic architecture with 4,096 neurosynaptic cores, totaling 1 million neurons and 256 million synapses [5]. The architecture employs a strictly digital design with binary synaptic weights and deterministic neuron dynamics. Each core operates independently with local memory and computation, communicating through an asynchronous Network-on-Chip (NoC) using Address-Event Representation (AER) protocol. TrueNorth's design prioritizes power efficiency through event-driven operation, consuming approximately 70mW at maximum activity. The chip operates at 1kHz biological real-time with deterministic timing, making it suitable for real-time edge applications. The architecture supports flexible neuron and synapse configurations, enabling implementation of various spiking neuron models including LIF, Izhikevich, and Hodgkin-Huxley dynamics [8].

C. BrainScaleS Analog Platform

BrainScaleS represents an alternative approach using analog circuit elements to emulate neuron and synapse dynamics [6]. The platform accelerates neural dynamics by factor 10,000 compared to biological real-time, enabling rapid network optimization and parameter exploration. The analog implementation achieves exceptional energy efficiency for synaptic operations, consuming approximately 20nW per synapse. However, analog neuromorphic platforms face challenges including parameter mismatch due to device variation and limited precision compared to digital implementations. BrainScaleS addresses these limitations through calibration mechanisms and hybrid analog-digital architectures that combine the energy efficiency of analog computation with digital control and communication [9].

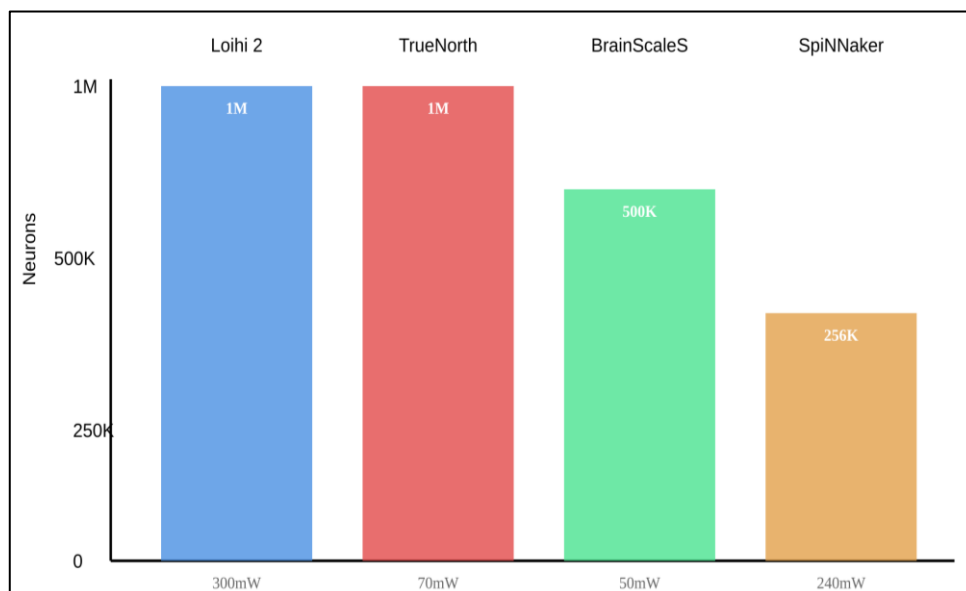


Fig. 1: Comparison of neuromorphic hardware architectures showing neuron counts, synaptic density, and power consumption.

III. SPIKING NEURAL NETWORK MODELS

A. Leaky Integrate-and-Fire Neurons

The Leaky Integrate-and-Fire (LIF) neuron model serves as the fundamental computational unit in most neuromorphic systems [2]. The membrane potential $V(t)$ evolves according to the differential equation:

$$\tau_m \frac{dV}{dt} = -(V - V_{rest}) + R_m I(t) \quad (1)$$

where τ_m represents the membrane time constant, V_{rest} is the resting potential, R_m denotes membrane resistance, and $I(t)$ represents input current. When $V(t)$ reaches the threshold V_{th} , the neuron emits a spike and resets to V_{reset} . The LIF model captures essential spiking dynamics while maintaining computational simplicity suitable for hardware implementation [10].

B. Spike-Timing-Dependent Plasticity

Spike-Timing-Dependent Plasticity (STDP) enables unsupervised learning in SNNs by modifying synaptic weights based on relative spike timing between pre-synaptic and post-synaptic neurons [7]. The weight change Δw follows an asymmetric temporal window:

$$\Delta w = A_+ \exp\left(-\frac{\Delta t}{\tau_+}\right), \text{ for } \Delta t > 0 \text{ (pre before post)} \quad (2)$$

$$\Delta w = -A_- \exp\left(\frac{\Delta t}{\tau_-}\right) \text{ for } \Delta t < 0 \text{ (post before pre)} \quad (3)$$

Where Δt represents the time difference between pre-synaptic and post-synaptic spikes, A_+/A_- control learning rates, and τ_+/τ_- determine temporal windows. This biologically inspired learning rule enables autonomous feature extraction and pattern recognition in neuromorphic systems [11].

C. Spike Encoding Mechanisms

Converting continuous sensor data into spike trains requires efficient encoding mechanisms [12]. Rate coding represents signal intensity through spike frequency, providing robustness but sacrificing temporal precision. Temporal coding encodes information in precise spike timing, enabling rapid inference but requiring accurate timing mechanisms. Latency coding represents stimulus intensity through time-to-first-spike, offering fast recognition with minimal spikes. Population coding distributes information across multiple neurons, providing noise robustness through ensemble averaging. Selection of appropriate encoding strategies depends on application requirements and available hardware resources.

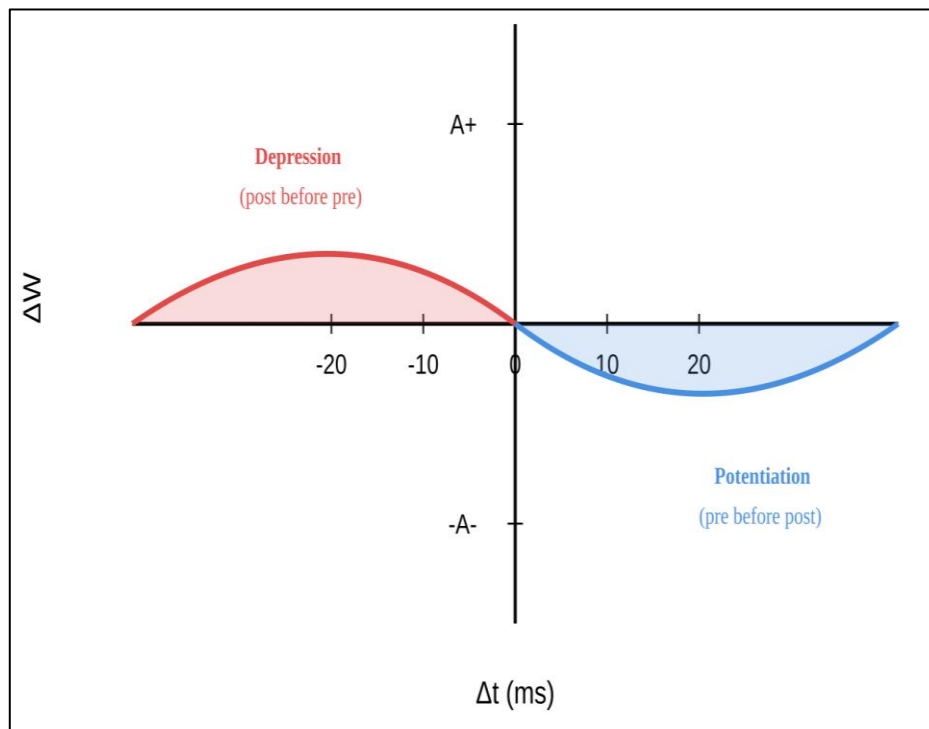


Fig. 2: Spike-timing-dependent plasticity learning window showing synaptic weight changes as a function of spike timing difference.

IV. PERFORMANCE AND ENERGY EFFICIENCY

A. Energy Consumption Analysis

Neuromorphic platforms demonstrate remarkable energy efficiency compared to conventional deep learning accelerators. Intel Loihi 2 achieves 0.26pJ per synaptic operation, representing 1000× improvement over GPU implementations at 260pJ per operation [3]. The energy efficiency stems from event-driven computation, where power consumption scales with network activity rather than peak computational capacity.

B. Inference Accuracy and Latency

While neuromorphic systems excel in energy efficiency, achieving competitive accuracy requires careful network design and training methodologies. Conversion-based approaches translate pre-trained ANNs to SNNs, achieving 98.4% accuracy on MNIST and 92.7% on CIFAR-10 with appropriate normalization and threshold balancing [13]. Direct SNN training using surrogate gradients attains comparable accuracy while maintaining temporal dynamics essential for neuromorphic hardware [14].

Inference latency in SNNs depends on spike propagation time and network depth. Shallow networks achieve sub-millisecond latency for simple classification tasks, while deep SNNs require 10-50 timesteps for convergence. Optimized spike encoding and early stopping mechanisms reduce average latency to 2-5ms for edge vision applications [15].

V. EDGE DEPLOYMENT STRATEGIES

A. Network Architecture Optimization

Deploying SNNs on resource-constrained edge devices requires architectural optimization to balance accuracy, latency, and energy consumption. Network pruning techniques remove redundant synaptic connections, reducing memory footprint and computational requirements while maintaining accuracy [16]. Weight quantization to 4-8 bits further decreases memory bandwidth without significant performance degradation. Sparse connectivity patterns inspired by biological cortical structure achieve 10-20× parameter reduction compared to fully-connected architectures.

B. Hybrid Computing Frameworks

Practical edge AI systems often combine neuromorphic accelerators with conventional processors in heterogeneous architectures [17]. Pre-processing stages execute on general-purpose cores, while SNNs handle feature extraction and classification on neuromorphic hardware. This hybrid approach leverages the strengths of each computing paradigm: conventional processors provide flexibility for complex control logic, while neuromorphic chips deliver energy-efficient inference for pattern recognition tasks.

C. Application Domains

Neuromorphic computing demonstrates particular promise in several edge AI domains. Event-based vision sensors combined with SNNs enable ultra-low-power object detection and tracking, consuming less than 1mW for continuous operation [18]. Keyword spotting applications achieve always-on voice activation with 200μW power consumption. Gesture recognition systems process temporal dynamics efficiently through recurrent SNN architectures. Predictive maintenance applications leverage STDP for online learning and adaptation to changing environmental conditions without cloud connectivity.

VI. CONCLUSION AND FUTURE DIRECTIONS

Neuromorphic computing architectures implementing spiking neural networks offer compelling advantages for energy-efficient AI inference at the edge. Current platforms demonstrate 100-1000× energy efficiency improvements compared to conventional deep learning accelerators while maintaining competitive accuracy for pattern recognition tasks. The event-driven computation paradigm aligns naturally with sparse, temporal data from edge sensors, enabling continuous processing with minimal power consumption.

Future research directions include developing scalable training algorithms that leverage neuromorphic hardware acceleration, establishing standardized benchmarks for comparing neuromorphic platforms, and creating software frameworks that abstract hardware-specific details. Heterogeneous integration of neuromorphic accelerators with conventional processors promises to unlock new application domains requiring both flexibility and extreme energy efficiency. As neuromorphic technology matures, we anticipate widespread deployment in battery-powered IoT devices, enabling intelligent edge processing without cloud dependency.

The convergence of advanced fabrication technologies, bio-inspired learning algorithms, and domain-specific architectures positions neuromorphic computing as a transformative approach for sustainable artificial

intelligence. Continued innovation in neuromorphic hardware and algorithms will enable intelligent edge devices capable of autonomous learning and inference while operating within severe power budgets essential for ubiquitous deployment.

REFERENCES

- [1] M. Davies et al., "Advancing neuromorphic computing with Loihi: A survey of results and outlook," *Proceedings of the IEEE*, vol. 109, no. 5, pp. 911-934, May 2021.
- [2] W. Maass, "Networks of spiking neurons: The third generation of neural network models," *Neural Networks*, vol. 10, no. 9, pp. 1659-1671, 1997.
- [3] S. K. Esser et al., "Convolutional networks for fast, energy-efficient neuromorphic computing," *Proceedings of the National Academy of Sciences*, vol. 113, no. 41, pp. 11441-11446, 2016.
- [4] M. Davies et al., "Loihi 2: A scalable neuromorphic processor," in *IEEE International Solid-State Circuits Conference (ISSCC)*, 2023, pp. 36-38.
- [5] P. A. Merolla et al., "A million spiking-neuron integrated circuit with a scalable communication network and interface," *Science*, vol. 345, no. 6197, pp. 668-673, Aug. 2014.
- [6] J. Schemmel et al., "Live demonstration: A scaled-down version of the BrainScaleS wafer-scale neuromorphic system," in *IEEE International Symposium on Circuits and Systems (ISCAS)*, 2012, p. 702.
- [7] G.-Q. Bi and M.-M. Poo, "Synaptic modifications in cultured hippocampal neurons: Dependence on spike timing, synaptic strength, and postsynaptic cell type," *Journal of Neuroscience*, vol. 18, no. 24, pp. 10464-10472, 1998.
- [8] S. Modha et al., "Cognitive computing," *Communications of the ACM*, vol. 54, no. 8, pp. 62-71, Aug. 2011.
- [9] S. Schmitt et al., "Neuromorphic hardware in the loop: Training a deep spiking network on the BrainScaleS wafer-scale system," in *International Joint Conference on Neural Networks (IJCNN)*, 2017, pp. 2227-2234.
- [10] M. Izhikevich, "Simple model of spiking neurons," *IEEE Transactions on Neural Networks*, vol. 14, no. 6, pp. 1569-1572, Nov. 2003.
- [11] S. Song, K. D. Miller, and L. F. Abbott, "Competitive Hebbian learning through spike-timing-dependent synaptic plasticity," *Nature Neuroscience*, vol. 3, no. 9, pp. 919-926, Sep. 2000.
- [12] Amir et al., "A low power, fully event-based gesture recognition system," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 7388-7397.
- [13] Rueckauer et al., "Conversion of continuous-valued deep networks to efficient event-driven networks for image classification," *Frontiers in Neuroscience*, vol. 11, p. 682, 2017.
- [14] E. O. Neftci, H. Mostafa, and F. Zenke, "Surrogate gradient learning in spiking neural networks," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 51-63, Nov. 2019.
- [15] P. U. Diehl and M. Cook, "Unsupervised learning of digit recognition using spike-timing-dependent plasticity," *Frontiers in Computational Neuroscience*, vol. 9, p. 99, Aug. 2015.
- [16] Y. Kim et al., "Spiking neural network using synaptic pruning and growth for classifying MNIST handwritten digits," in *IEEE International Conference on Big Data and Smart Computing (BigComp)*, 2020, pp. 31-36.
- [17] Frenkel, M. Lefebvre, J.-D. Legat, and D. Bol, "A 0.086-mm² 12.7-pJ/SOP 64k-synapse 256-neuron online-learning digital spiking neuromorphic processor in 28-nm CMOS," *IEEE Transactions on Biomedical Circuits and Systems*, vol. 13, no. 1, pp. 145-158, Feb. 2019.
- [18] Gallego et al., "Event-based vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 1, pp. 154-180, Jan. 2022.