

# Scaling the Spike: Unlocking Neuromorphic Computing at Brain-Scale

Satyajit Deokar

## 1. Toughest feature to scale – and why

Among the highlighted features—**distributed hierarchy**, **sparsity**, **neuronal scalability**, **plasticity**, and **interconnectivity**—**neuronal scalability** stands out as the most challenging. It requires supporting hundreds of millions (or billions) of neurons across chips and racks while maintaining real-time dynamics and low energy use ([gvern.net](http://gvern.net)). The difficulties lie in:

- **Inter-chip communication bandwidth** needed to preserve spiking timing,
- **Synchronizing event-driven systems** without central clocks,
- **Efficient memory integration** to support synaptic plasticity at scale.

Successfully overcoming this will unlock:

- Real-time simulations of full-brain-scale networks,
- Energy-efficient, on-device intelligence for robotics, IoT, and edge AI,
- Solving NP-hard and graph-based problems natively on neuromorphic hardware.

## 2. The “AlexNet moment” for neuromorphic computing

Analogous to deep learning’s surge tied to AlexNet and GPUs, neuromorphic computing is waiting on breakthroughs in:

- **Programmable, memristor-based synaptic arrays** co-integrated with CMOS,
- **Scalable, asynchronous interconnect architectures** (e.g. hierarchical wafer-scale switches),
- **Event-driven learning algorithms** that rival backpropagation in efficiency.

A plausible trigger could be a system combining memristive crossbar arrays (for in-memory synaptic updates), local plasticity rules, and a deep spiking-learning framework—leading to breakthroughs in:

- Ultra-low-power vision,
- Real-time sensory integration (audio-vision),
- Always-on robotics requiring adaptive on-device intelligence (e.g., drones, prosthetics).

## 3. Proposal for hardware–software framework interoperability

To close the implementation–software gap:

1. **Define a common intermediate representation (IR)**—a spiking graph model with layers, synapse types, neuron models, and routing constraints.
2. **Develop compiler backends** for existing chips (Loihi, TrueNorth, SpiNNaker) to translate IR to hardware-specific deployments.
3. **Extend frameworks like PyTorch or TensorFlow** with a “Neuromorphic mode,” offering constructs for spiking layers and event-based training.
4. **Push open-source connectors** (e.g. ONNX-neuromorphic) that target both hardware simulators and physical chips.
5. **Launch community challenges** (e.g. SNN-ImageNet, event-based object detection) to validate interoperability across platforms.
- 6.

# Scaling the Spike: Unlocking Neuromorphic Computing at Brain-Scale

## Satyajit Deokar

### 4. Unique benchmarking metrics beyond accuracy

Neuromorphic systems demand metrics reflecting their core strengths:

Metric	Motivation
Energy per inference + learning	Captures ultra-low power event-driven gains
Spike latency	Measures time from input event to output, essential for low-latency contexts
Sparsity rate	Ratio of active neurons/synapses—shows efficiency
Plasticity/synaptic updates per second	Evaluates learning capability
Scalability score	Efficiency normalized across simulated neuron counts
Fault tolerance	Performance degradation under simulated failures

Standardization via **Tiered Benchmarks**:

- Tier 1: Basic tasks (e.g. MNIST, N-back recognition),
- Tier 2: Event-based vision/robotics,
- Tier 3: Large spiking network applications (e.g. graph problems),

with mandatory reporting across all metrics to allow fair comparison across architectures.

### 5. Why emerging memories matter

Emerging memories like **memristors** and **phase-change memory (PCM)** enable **compute-in-memory** and **in-situ plasticity**, sidestepping the bottleneck of traditional von Neumann designs. Specific promising avenues:

- **Analog multi-bit synaptic weights**: enabling dense plastic arrays,
- **On-chip STDP mechanisms**: supporting local learning,
- **3D-stacked arrays** co-located with neuron circuits for tight integration,
- **Stochastic memristive behavior** for inherently probabilistic spiking models (e.g. Boltzmann-like inference).

Key research directions:

- **Variability-aware learning**, making use of device imperfections,
- **Hybrid CMOS-memristor neurosynaptic cores** with local event-driven plasticity,
- **Mapping deep SNNs onto memristor hardware**, balancing precision and sparsity,
- **Exploring PCM for homeostatic and metaplastic mechanisms**, enabling self-organizing neuromorphic systems.

## **Scaling the Spike: Unlocking Neuromorphic Computing at Brain-Scale**

**Satyajit Deokar**

### **Final Take**

Neuromorphic computing at scale is at a pivotal inflection point—poised for breakthroughs if we can engineer scalable synaptic fabrics, co-design software–hardware ecosystems, and define benchmarks that capture real-world dynamics. With integrated memories and event-driven plasticity, neuromorphic systems could leap beyond von Neumann limits to deliver energy-efficient, adaptive, and trustworthy intelligence in edge and networked environments.