

Merging insights from artificial and biological neural networks

Does neuromorphic edge intelligence need spikes?

Charlotte Frenkel (<u>c.frenkel@tudelft.nl</u>)

Assistant Professor Dept. of Microelectronics, Delft University of Technology

NICE 2025, March 27th

Outline

① What are the key synergies between the bottom-up and top-down approaches?

2 How can we exploit these synergies for novel spike-based engineering solutions?

Synaptic plasticity rules – Neuroscience as the starting point



Frenkel, NICE 2025

3

Neural network training – Bio-plausibility as the end goal Synergy with hardware: latency, memory access patterns



Frenkel, NICE 2025

[Lillicrap, Nat. Comms., 2016] [Nokland, NeurIPS, 2016] [Frenkel & Lefebvre, Front. Neur., 2020]

AI algorithms

Neuroscience

HW efficiency and bio-plausibility are often two sides of the same coin!



On our way to neuromorphic intelligence – Bottom-up or top-down?



Outline

① What are the key synergies between the bottom-up and top-down approaches?

2 How can we exploit these synergies for novel spike-based engineering solutions?

Let's use a 4-step recipe!

 Neuromorphic intelligence:

 2 should be fed by 1





Different users, environments, task requirements

More training data before deployment?

Issues: cost, robustness, flexibility

Data exchange with the cloud?

Issues: power budget, privacy



Why is on-chip learning over second-long timescales difficult? Let's solve a yet unsolved engineering challenge!



- Unrolling in time: very deep network (current learning ICs for static stimuli: ≤3 layers)
- Intractable memory/latency requirements
- No end-to-end on-chip solution to date

<u>Key challenge</u>: On-chip learning over long timescales while keeping a fine-grained temporal resolution

2) Select the (ML-informed) starting point

From BPTT to biologically plausible training



3) Use-case-driven feature set selection

Neuron model selection... driven by the application requirements!



4) Enforce space and time locality

Key steps to minimize memory requirements



Stochastic weight updates allow reducing weight resolution to 8 bits

The ReckOn neuromorphic chip – Microphotograph and summary



					_
Technology	28nm FDSOI CMOS				
Core size	0.67 x 0.67 mm ² 0.45mm ²				
Die size	0.93 x 0.93 mm ²				
SRAM		138kB	+ 0	kB ext. D	RAM!
Network	Spiking RNN				
Training timespan	Max. 32k steps				
					-



The ReckOn neuromorphic chip – Key advantage of using spikes



- Event-driven / sparsity-aware computation
- Sensor-agnostic raw-data processing
- Task-agnostic processing and learning



Neuromorphic spiking sensor

[Frenkel, *ISSCC*, 2022]

The ReckOn neuromorphic chip – Benchmarking



Outline

① What are the key synergies between the bottom-up and top-down approaches?

2 How can we exploit these synergies for novel spike-based engineering solutions?

Let's now look into spikes for low-latency applications!



Chauvaux

[N. Chauvaux et al., ISCAS'25 (accepted)] Extension preprints coming soon.

Spiking neural networks for low-latency event-driven computation



have a memory overhead penalty.

A good case for bringing compute close to memory!

Compute in memory (CIM) to the rescue of spiking neural networks







🖹 Sparse inputs?

X Non-linearity, noise, mismatch



- Low parallelism... \mathbf{x}
 - ...but efficient event-driven
 - addition and accumulation!
 - Robustness

Current digital CIM approaches for SNNs have a flexibility issue



Unlocking dataflow and resolution flexibility in digital CIM for SNNs



Solution:



👽 Utilization can be maximized



Unlocking dataflow and resolution flexibility in digital CIM for SNNs





Weight

stationarity

or

Output

stationarity

FlexSpIM – A flexible spiking in-memory macro





TSMC 40nm, proven in silico!



FlexSpIM can be tailored to the use case and target specifications!

Outline

① What are the key synergies between the bottom-up and top-down approaches?

2 How can we exploit these synergies for novel spike-based engineering solutions?

Let's now look into spikes for low-latency applications!



[Y. Yang*, A. Kneip*, C. Frenkel, Trans. CAS AI, 2025] Open-source: github.com/cogsys-tudelft/evgnn

Graphs – A better representation for event-based data?



From static graphs to dynamic graphs



<u>Challenge</u>: GNNs focus on static graphs, how we make them work for event-based data?

1) Graph building

2) GNN execution

N layers = information can propagate through N hops

For each event, only the N-hop neighborhood needs to be updated

AEGNN [Schaefer, CVPR, 2022]

AEGNN promises low-latency execution through local processing!



Directed dynamic graphs for low-latency hardware acceleration



<u>Challenge</u>: GNNs focus on static graphs, how we make them work for event-based data?

Graph building
 GNN execution
 AEGNN [Schaefer, CVPR, 2022]

But information flows from the future to the past...?

Let's use directed graphs! HUGNet [Dalgaty, CVPR-W, 2023]

Core insight

Now, data is *really* exchanged locally, in an event-driven fashion!

Core hardware contributions

- 1) Edge-free graph storage!
- 2) Neighborhood search decoupled in space-time
- 3) Parallel layer execution
- 4) No HW yet let's harvest OoMs with a simple design!

EvGNN – The first GNN-based hardware accelerator for event-based vision

Latency per event of 16µs (N-CARS, 87.8%) in a KV260 edge FPGA platform!

First hardware aligning with the temporal resolution of event-based cameras on a real-world benchmark. No custom silicon needed!

[Frenkel, ISSCC'22]

Yes, for sensor- and task-agnostic learning

[Chauvaux, ISCAS'25]

[Yang and Kneip, Trans. CASAI'25]

Yes, for low-latency event-based processing

The Cognitive Sensor Nodes and Systems (CogSys) Team

We bridge the bottom-up (bio-inspired) and top-down (engineering-driven) NeuroAI design approaches toward multi-scale, decentralized intelligence.

Funding acknowledgements

Frenkel, NICE 2025