

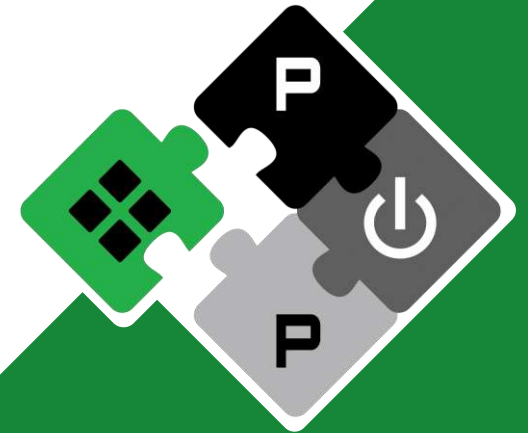
Neuromorphic Computing on PULP

Accelerating Spiking Neural Networks with digital Hardware

Alfio Di Mauro adimauro@iis.ee.ethz.ch

PULP Platform

Open Source Hardware, the way it should be!



@pulp_platform 

pulp-platform.org 

youtube.com/pulp_platform 

What is Neuromorphic Computing?



Artificial Neural Network

- Dimensions
 - Height
 - Width
 - Channels
- Non-Linear Activation
 - ReLu
 - Leaky-ReLu
 - Tanh
 - ...

Spiking Neural Network

- Dimensions:
 - Height
 - Width
 - Channels
 - **Time**
- **Neuromorphic Activation**
 - **Leaky Integrate & Fire**
 - **Hodgkin–Huxley**
 - ...

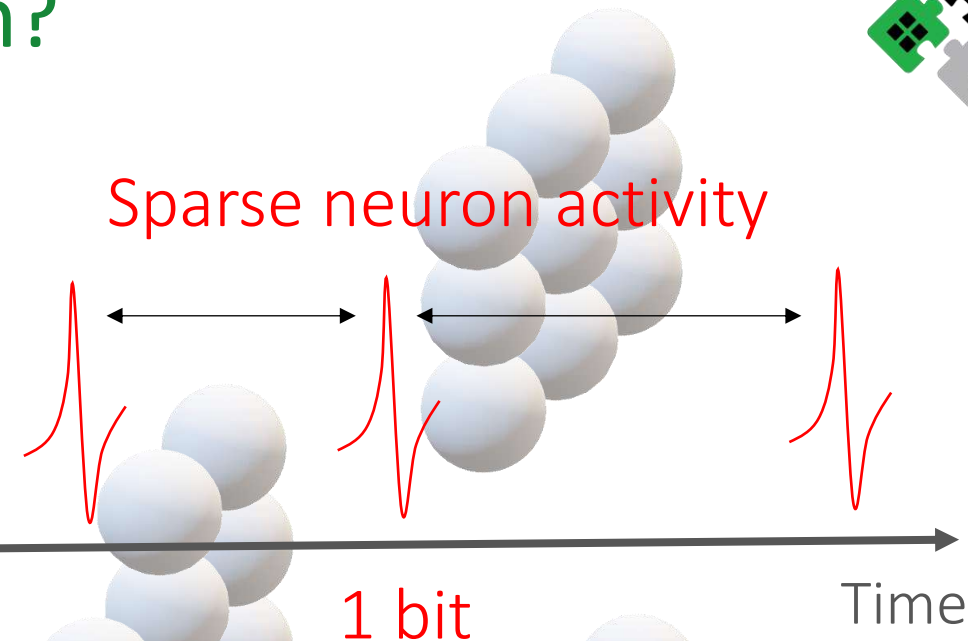
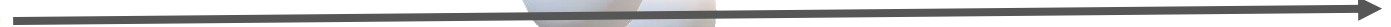


What is so special about SNNs then?



Leaky Integrate & Fire neuron

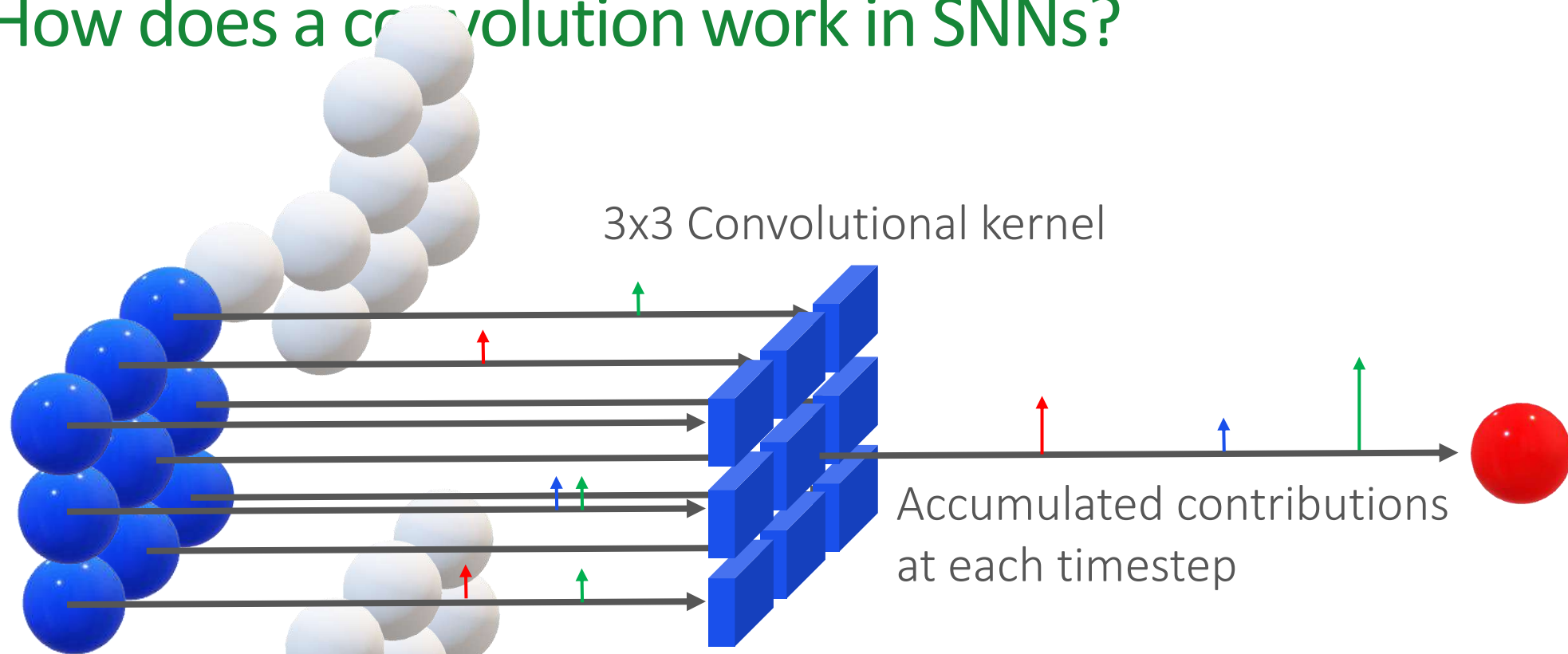
- Low complexity
- Temporal dynamics



Population of 6x6x1 LIF neurons



How does a convolution work in SNNs?



Binary tensor signalling which neuron is active at each timestep

- **Perform an addition only if a spike is present**
- **No multiplications**



Sparse data representation & Execution model



Event memory layout as a stream of
List of coordinate (COO) + time

OP	Time	Ch	Y	X
----	------	----	---	---

```
For k(c) in range(0, c_o):  
    load_sne(k(c), w0) # pre-load the weights
```

Output Channels (SW managed)

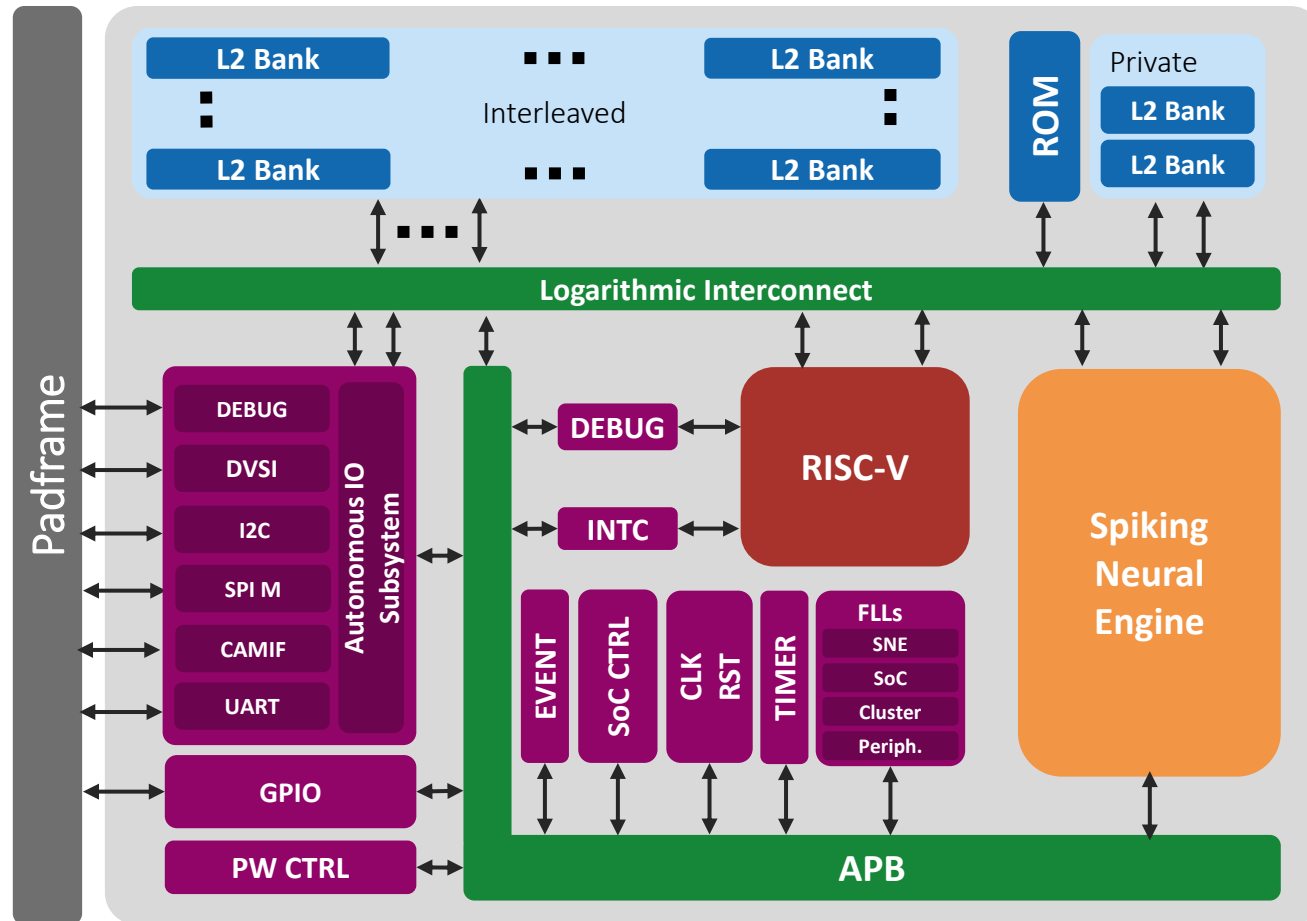
```
For evt_i in events_in: # fetch a new event from the memory  
    if evt_i is TIME_EVT: # check if it's a time event  
        t = evt_i # store the time  
        Events_out.push_time_evt(t) # propagate time on the output stream  
        Break # return from this iteration  
    Else:  
        c_i, e_x, e_y = get_spatial_coordinates(evt_i) # pre-process the event  
        # perform the output neuron update  
        For i in [0,7]: # span horizontal input event neighborhood  
            For j in [0,7]: # span vertical input event neighborhood  
                w0_ij = extract_weight(i,j,c_i,e_x,e_y,w0) # extract the respective weight for output ch 0  
                Events_out.push_evt(i,j, LIF_dynamic(w0_ij,t)) # write output events from ch 0 to the memory
```

Output neuron update

SNE execution



Can we accelerate SNNs on PULP?



HWPE Architecture

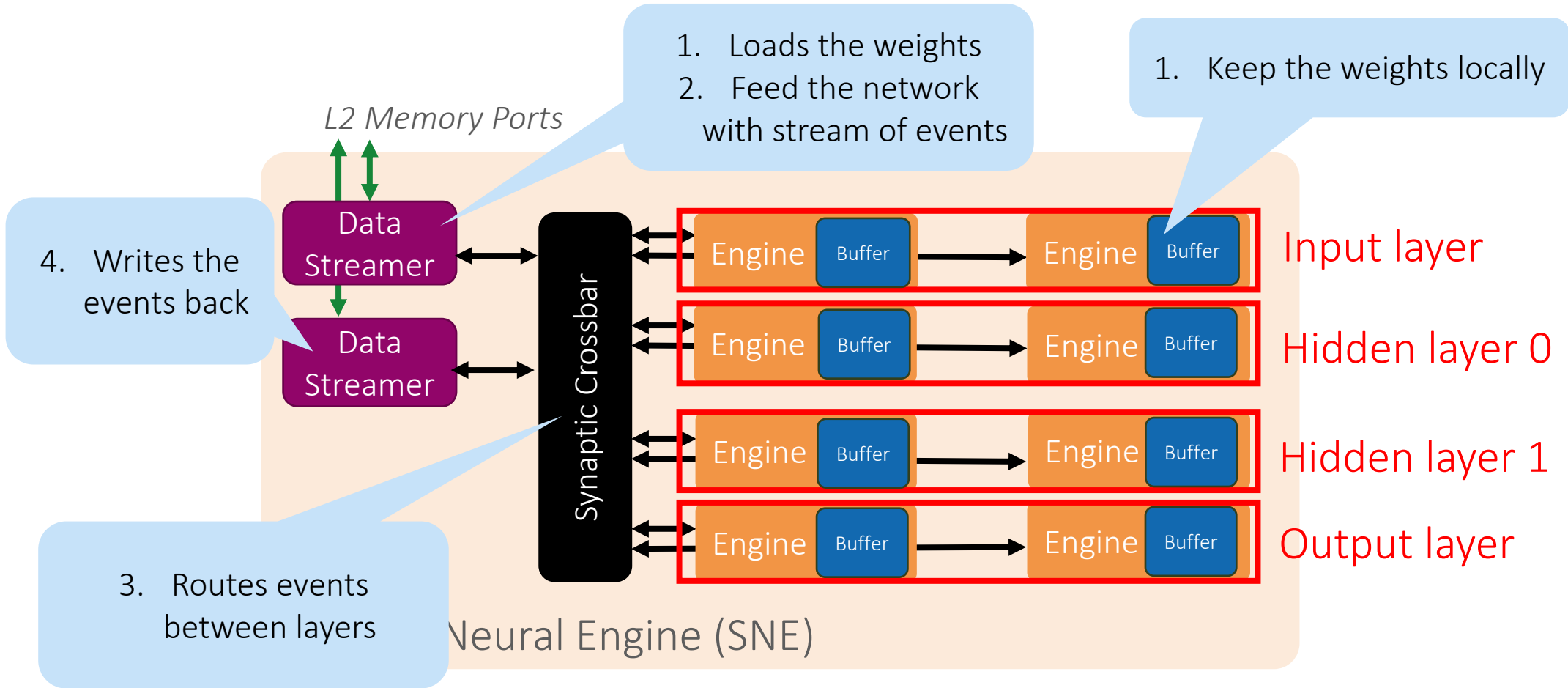
- ❑ High bandwidth to the L2 Memory/Peripherals
- ❑ Close interaction with a GP core

Challenges:

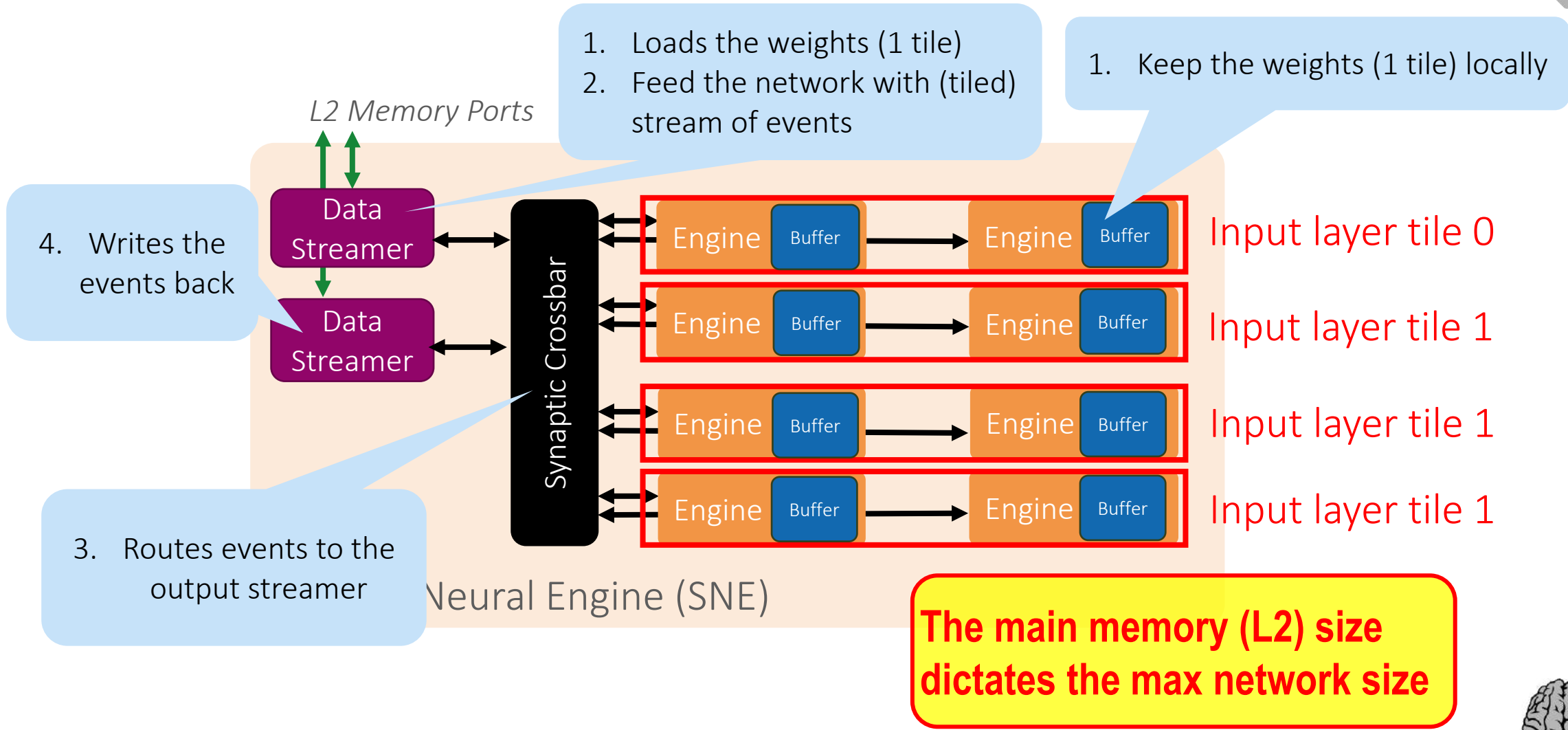
- ❑ Deal with data streams
- ❑ Exploit data sparsity
- ❑ Preserve energy-proportionality



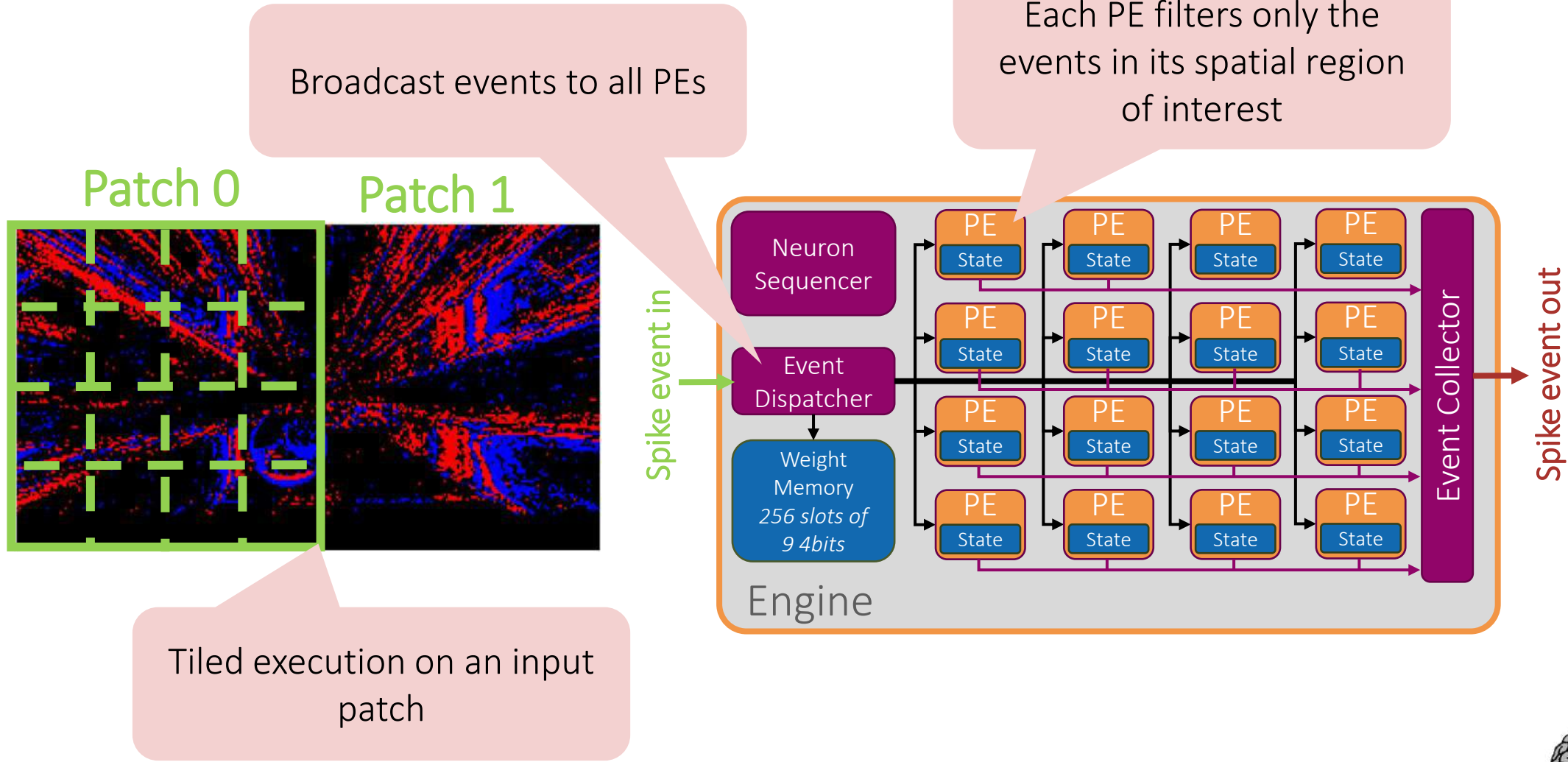
SNE is a data flow architecture



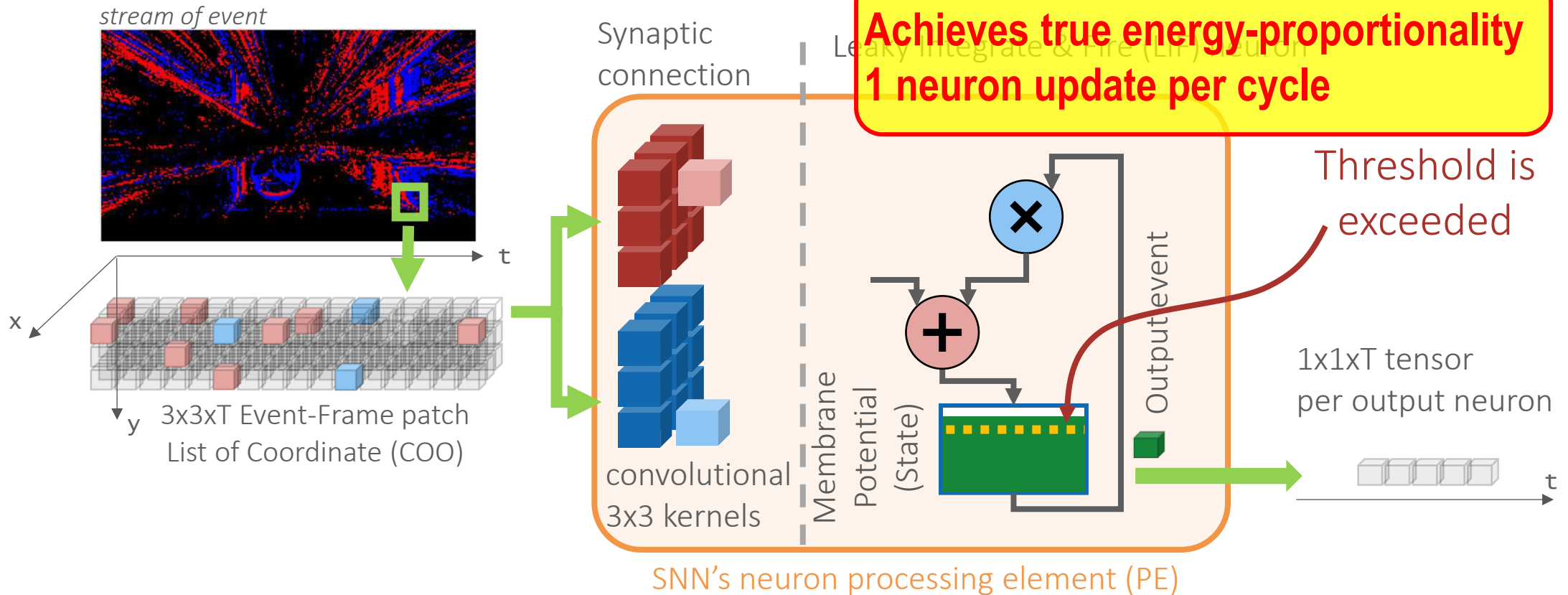
What if the network is too big?



Engine's Architecture



Event consumption, and output spikes generation



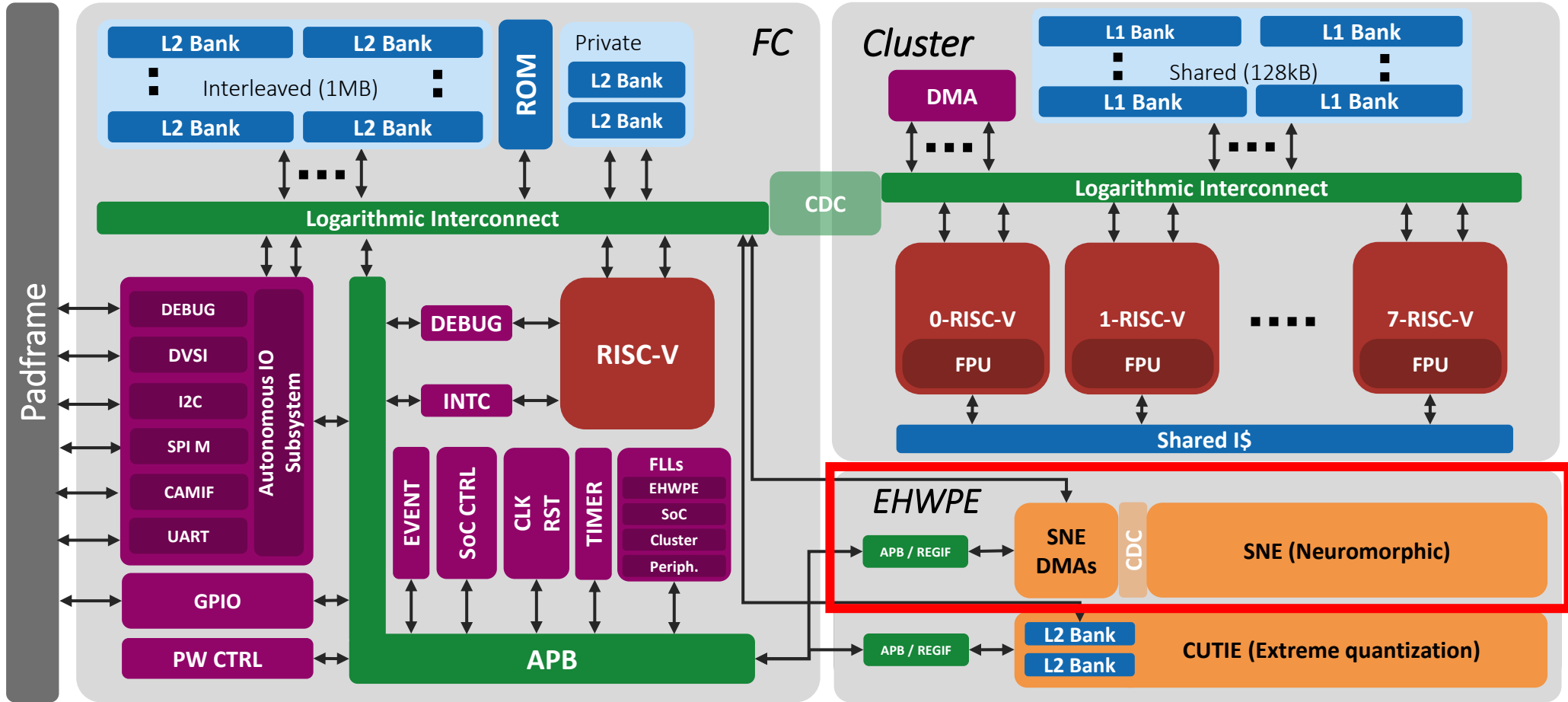
SNN's neuron processing element (PE)

A more complex dynamic than conventional DNNs neurons:

- Membrane Potential Accumulation/Activation **1 SynAcc = 1 4b-ADD + 1 8b-COMPARE**
- Membrane Potential decay **1 SynDec = (1 8b-MUL) + (1 8b-MUL + 1 8b-ADD)**

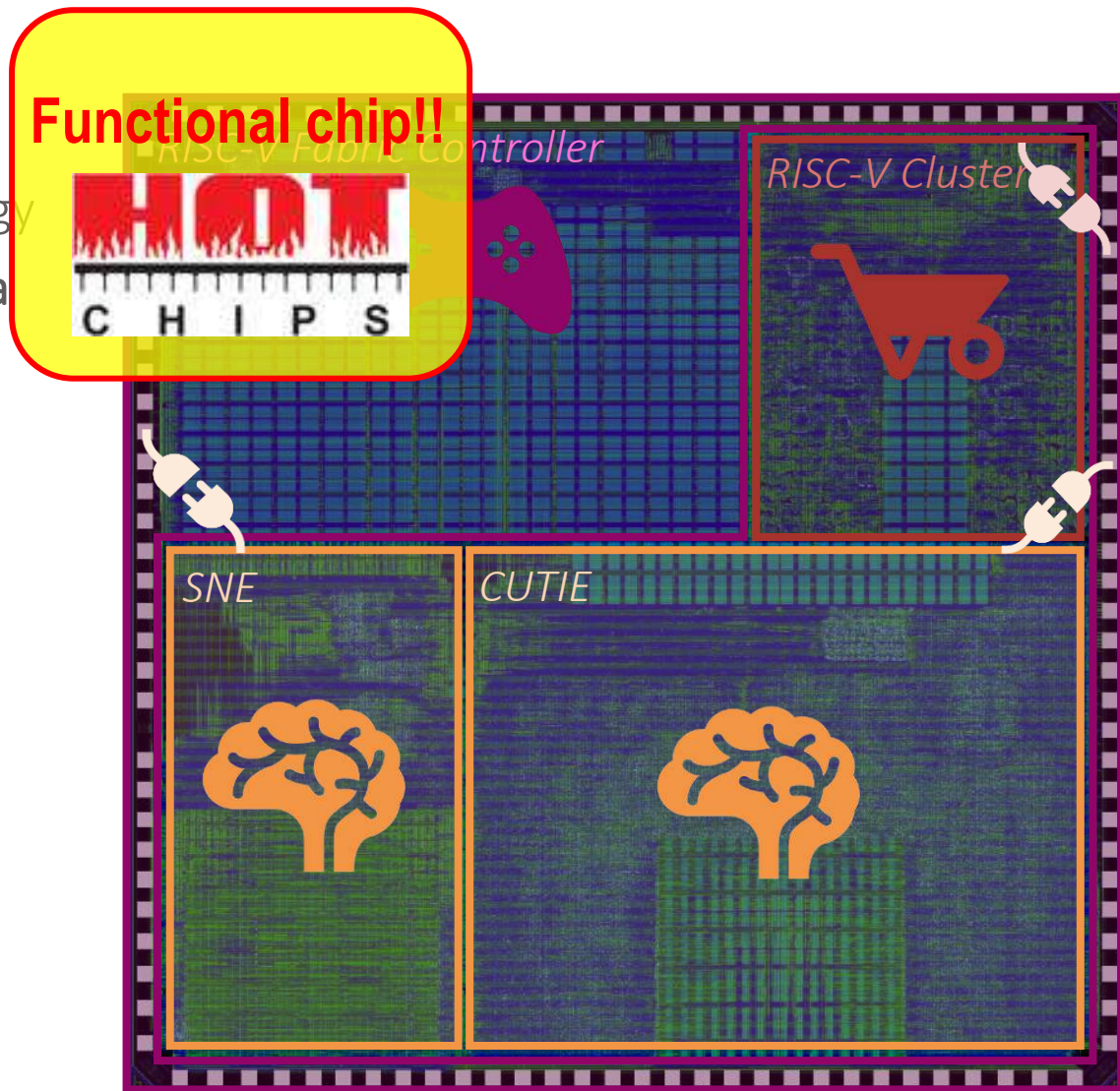


Alright we go to silicon: The Kraken Chip

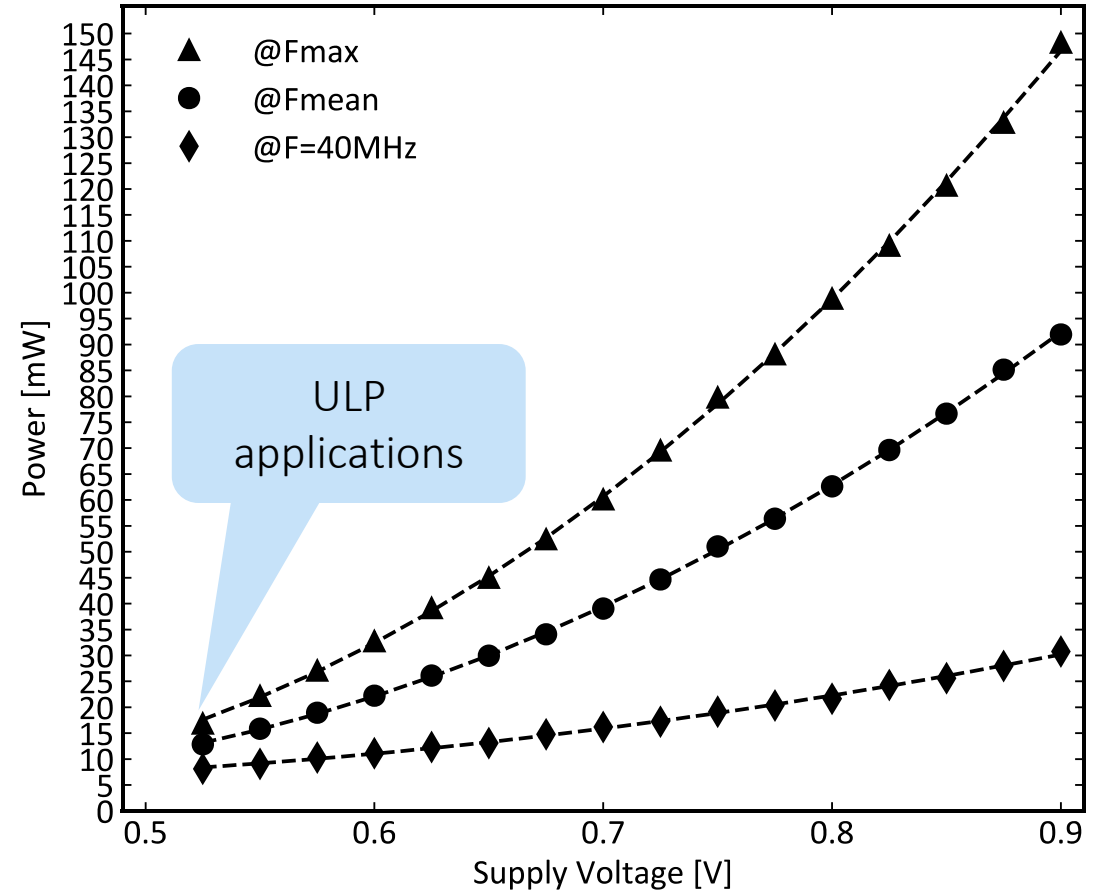
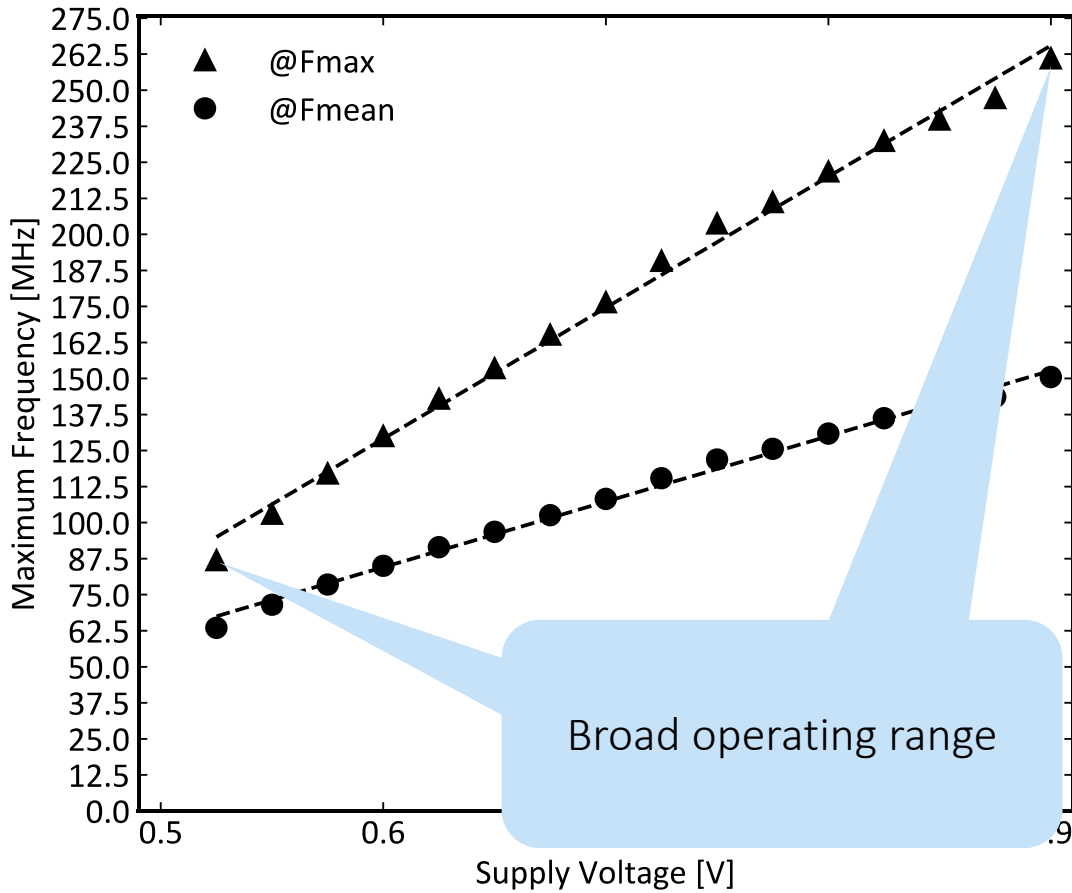


Silicon implementation

- GlobalFoundries 22nm FDX technology
- QFN88 chip package, 9mm² chip area
- 0.5V to 0.9V operating voltage
- Cluster Max Freq: **370MHz**
- CUTIE Max Freq: **140MHz**
- SNE Max Freq: **220MHz**
- Independent clock/power domain:
 - **RISC-V Cluster**
 - **SNE**
 - **CUTIE**



Power and Performance

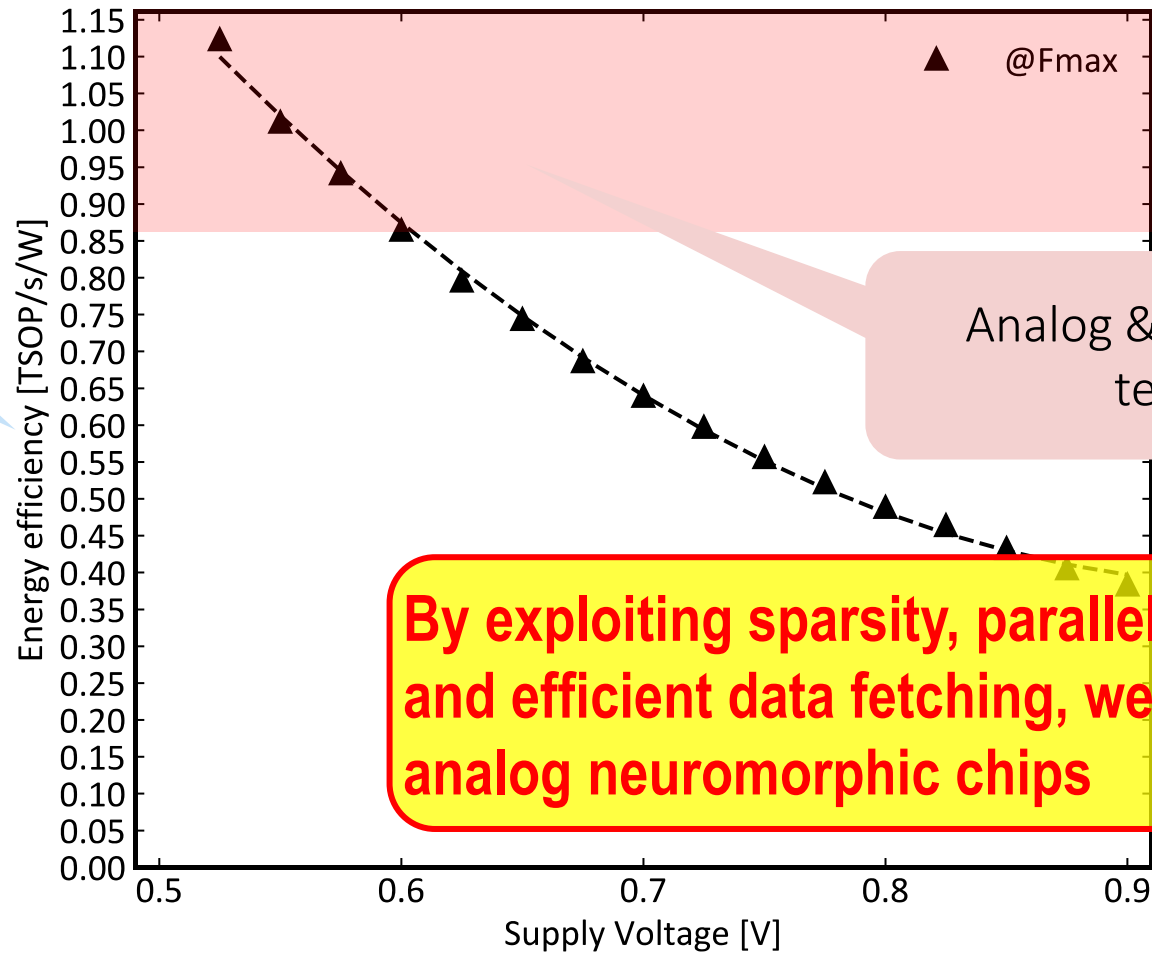


Yes, we can accelerate SNNs efficiently on PULP!



It's constant in SNE
Regardless of the
workload

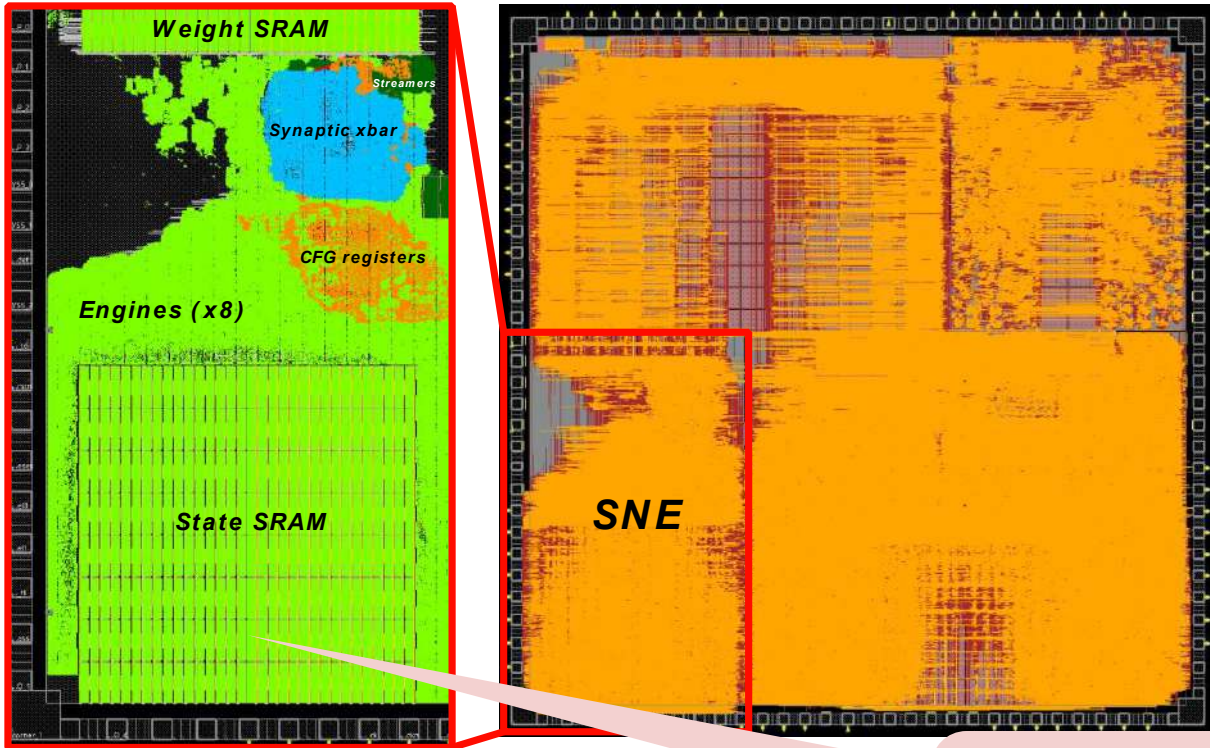
24 Synaptic Acc/cycle
522 Synaptic Dec/cycle



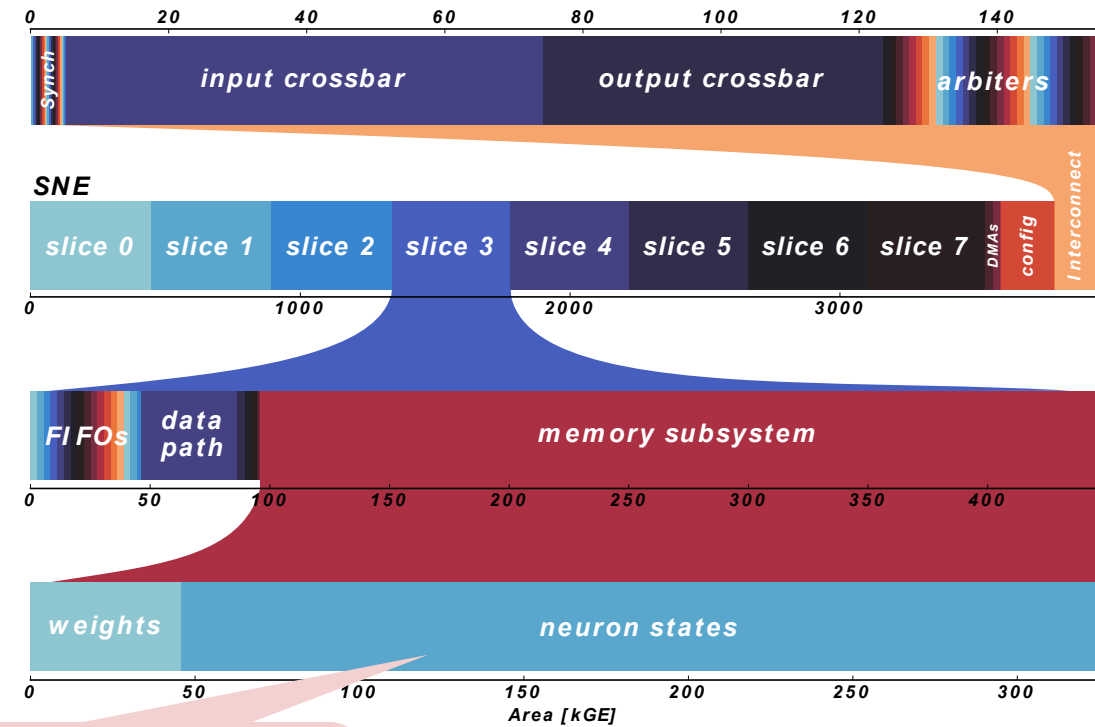
Area Breakdown



Kraken Layout



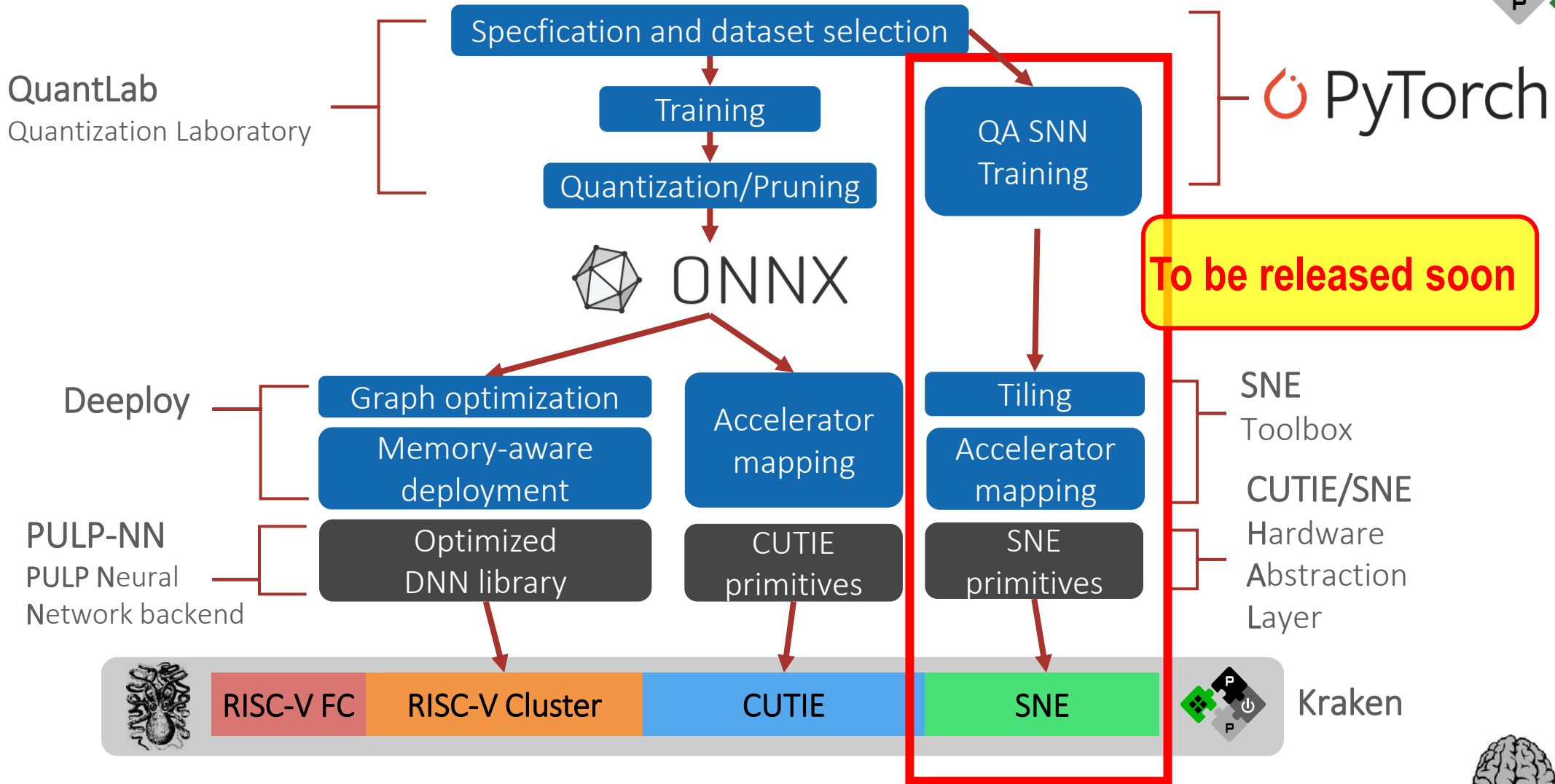
SNE Area (Kraken implem.)



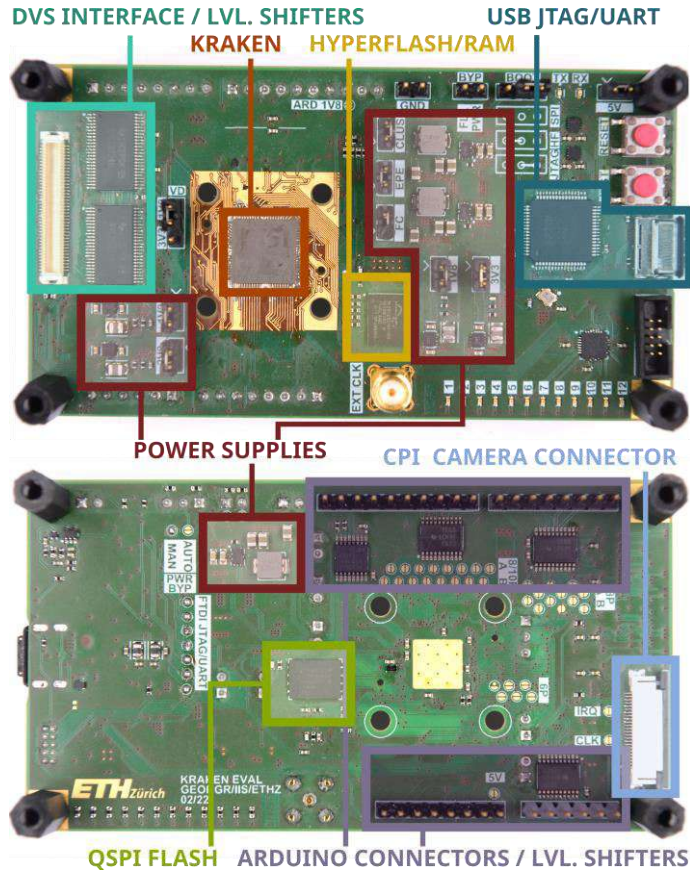
Yes, a lot of memories



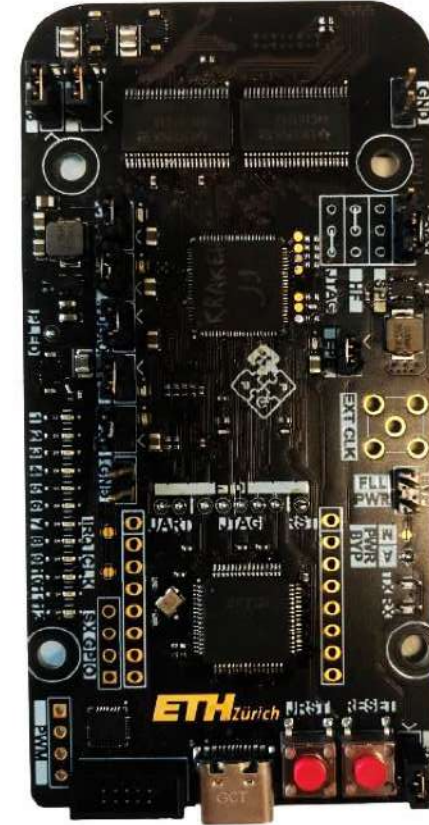
How do I train an SNNs for Kraken?



Release the Kraken



Kraken Eval Board v0



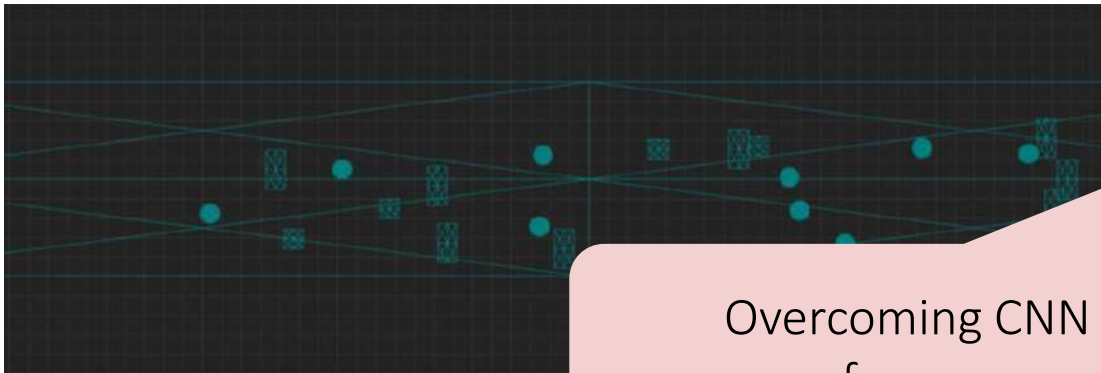
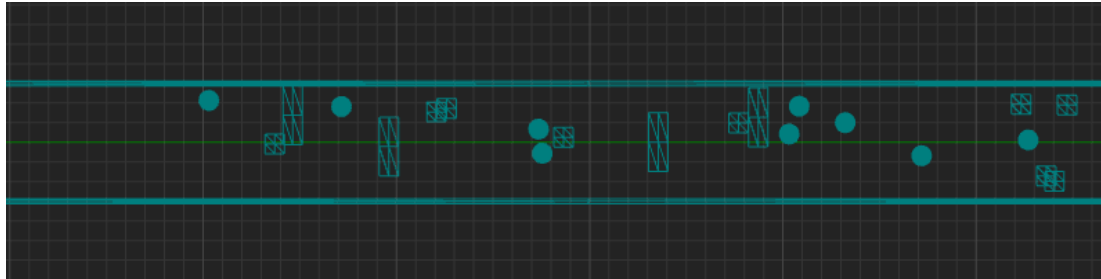
Kraken Eval Board v1



Applications: Reinforcement learning, Drone Navigation



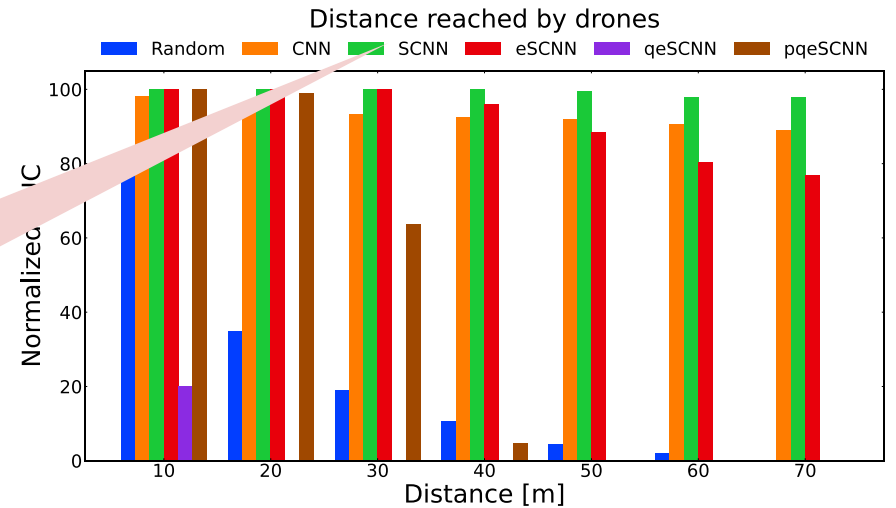
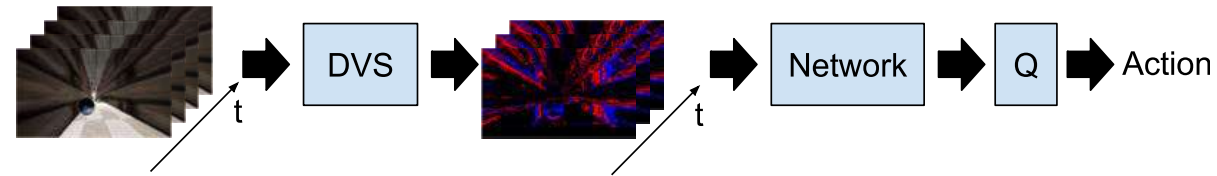
Corridor with obstacles



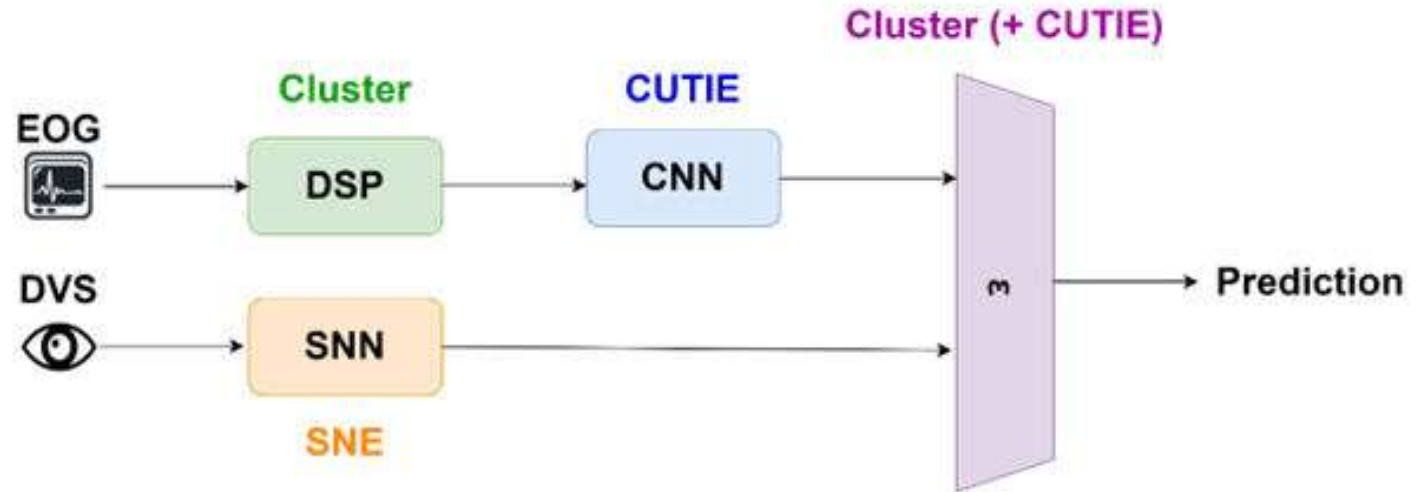
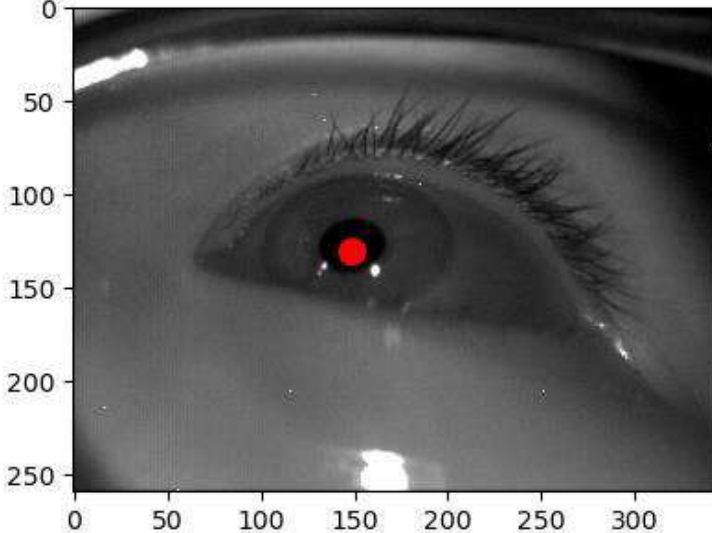
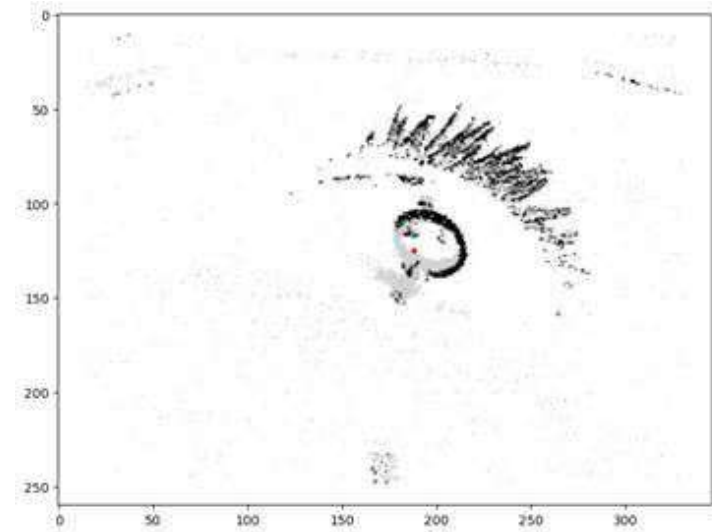
10 m

Overcoming CNN performance

Video processing pipeline (Kraken)



WIP Application: Eye tracking with fused modalities



Robust eye gaze prediction:

- Selective field of view Rendering in AR applications
- Human-Machine Interfaces



Conclusion

SNNs acceleration challenges:

- Exploit the sparse nature of the networks
- Perform computation in an energy-proportional way

Key take-away message:

- SNNs are not very different from Conventional ANNs
- Analog and Mixed-signal is not the only way for neuromorphic computing to be efficient

[Kraken: A Direct Event/Frame-Based Multi-sensor Fusion SoC for Ultra-Efficient Visual Processing in Nano-UAVs](https://www.research-collection.ethz.ch/handle/20.500.11850/565105)

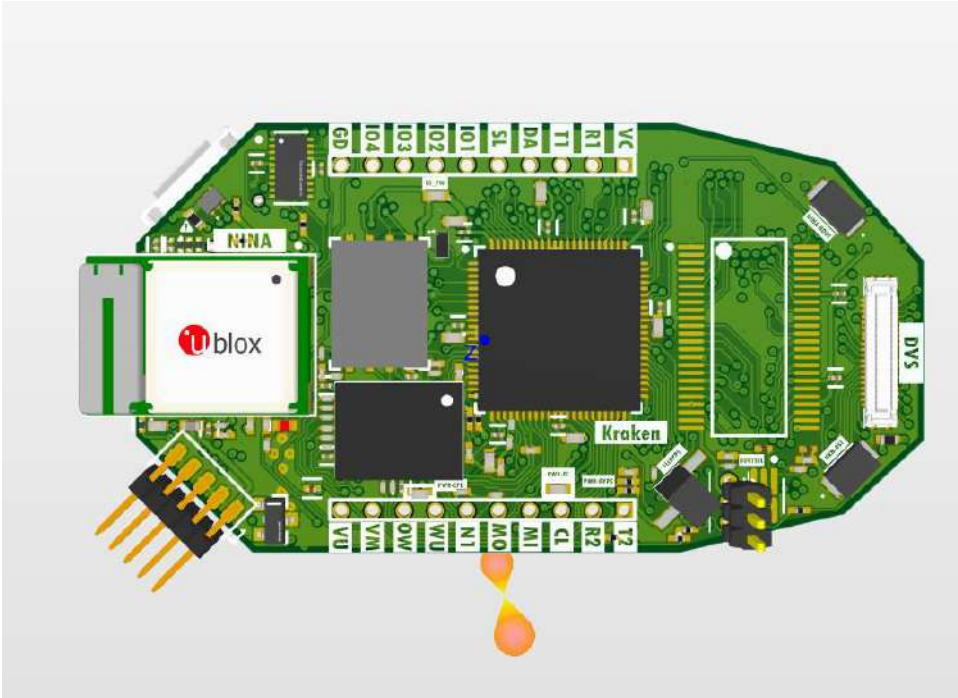
<https://www.research-collection.ethz.ch/handle/20.500.11850/565105>



[pulp-platform/sne \(github.com\)](https://github.com/pulp-platform/sne)



What's next? Kraken goes on a drone



Use Kraken to pilot a nano-drone in an autonomous drone race!!

- Event-based Camera
- Frame-based Camera
- SNNs (SNE)
- DNNs (CUTIE + Cluster)

