# AIGOR: Novel Neuromorphic Computing Architectures with a Multi-Node FPGA system

INFN Sezioni di Roma, Roma 2, APE Lab

AIGOR is a neuromorphic computer architecture prototype built on a multi-node FPGA (Field Programmable Gate Array) system. It's main objective is to serve as an accelerator platform for efficiently executing Spiking Neural Networks on edge devices. Brain-inspired Spiking Neural Networks represent a promising frontier in computational models, offering potential advantages over traditional computing paradigms in terms of energy efficiency, temporal information processing, and adaptability to dynamic data. This can benefit numerous applications, such as real-time signal processing and pattern recognition in resource-constrained environments. Neuromorphic computing is an approach to hardware architecture design to efficiently implement these biologically-inspired networks, balancing biological plausibility against computational efficiency. Leveraging on a proprietary framework for flexible, low latency communication we aim to deploy our architecture prototype on a multi-FPGA system, adopting a software-hardware codesign workflow that relies on the High Level Synthesis (HLS) programming paradigm for relatively fast and simple translation from a high level simulator of the architecture to the hardware design.

#### **Neuromorphic Computing**

#### **Spiking Neural Networks (SNNs):**

Artificial neural networks that communicate through "spikes" or pulses, similarly to biological neurons

#### Neuromorphic computing

A brain-inspired approach to hardware and algorithm design that efficiently realizes Spiking Neural Networks

**Brain:** processing-storage closely tied,

highly parallel 3D stacked architecture

High energy efficiency,

real-time processing,

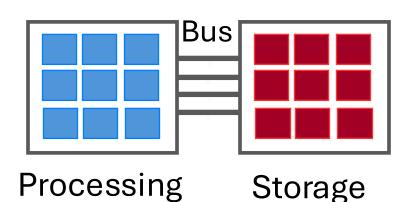
incremental learning

#### Why mimic the brain?

■ In-Memory Computing. Co-location of processing and storage

#### Traditional Von Neuman computing: The cost of data movement

Operation	Energy consumption
Addition of data	1x
Access data (onchip cache)	60x
Access data (offchip RAM)	3500x



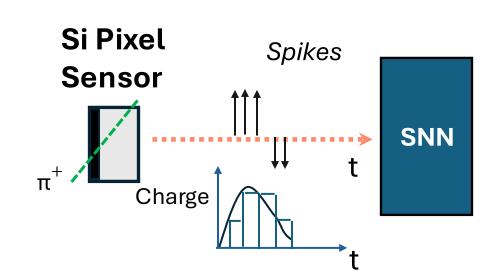
Processing

Sze et al, IEEE Custom Integrated Circuits Conference (2017)

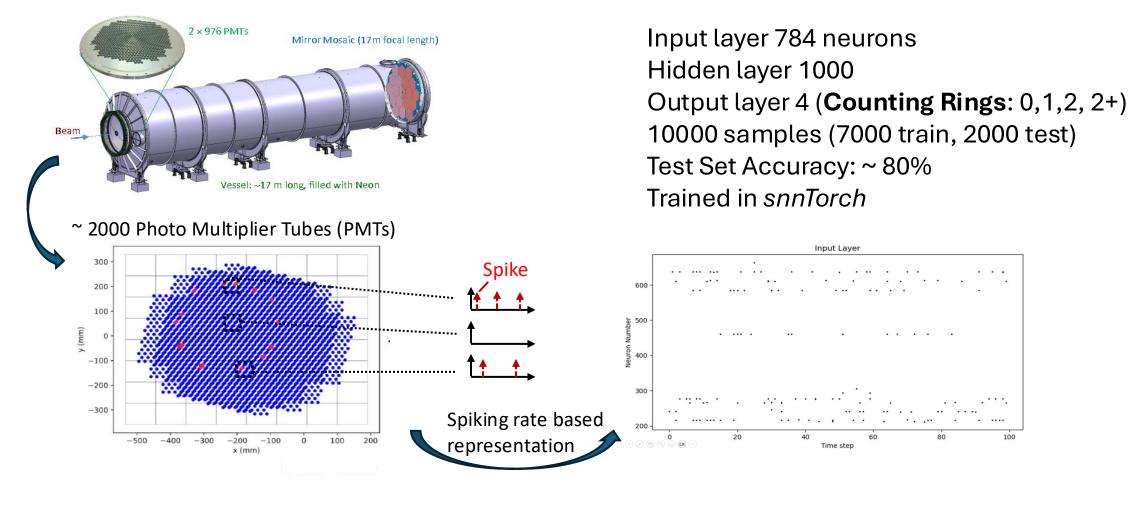
- Sparse, distributed information encoding through spikes
- Asynchronous, event-based computation relying on spikes time
- Local learning rules
  - No costly back-propagation

# **Neuromorphic Computing - Applications**

- Any input can be encoded as spikes
- Particularly effective when input is already sparse, event-based or is acquired as a time sequence
- Identify interesting applications in physics (particle physics sensors, anomaly detection, event cameras, time series data ...)



■ **Test-case**: NA62 Ring Imaging Cherenkov (RICH) detector

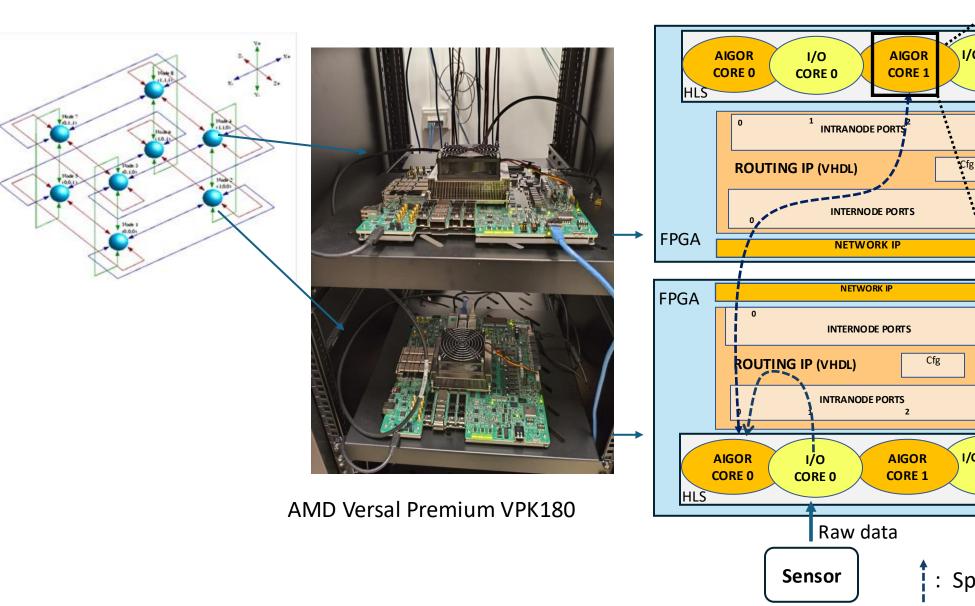


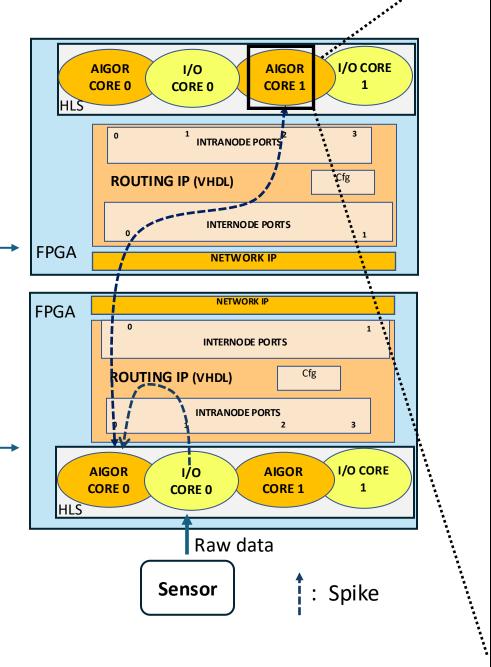
Ideal for edge Computing, on-sensor processing, streaming applications

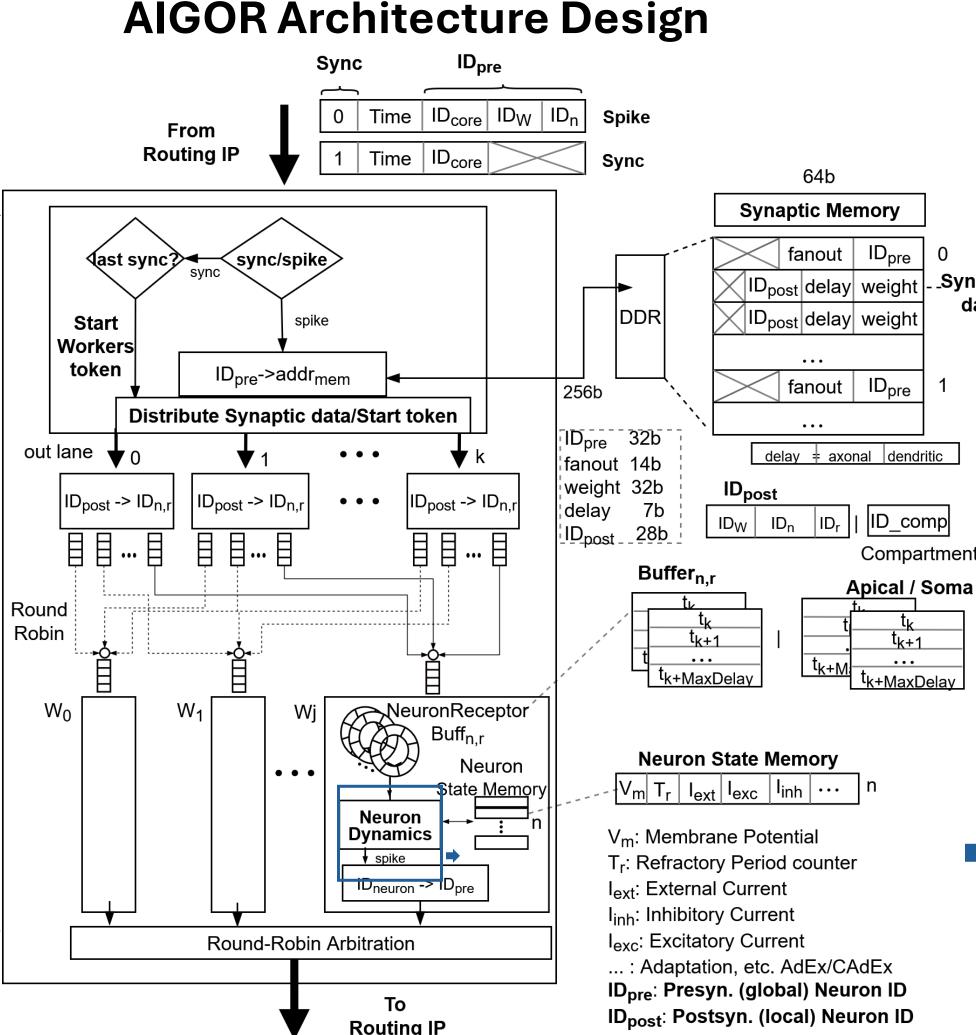
### Multi-FPGA architecture prototype

Direct network of interconnected FPGAs. 3D Torus.

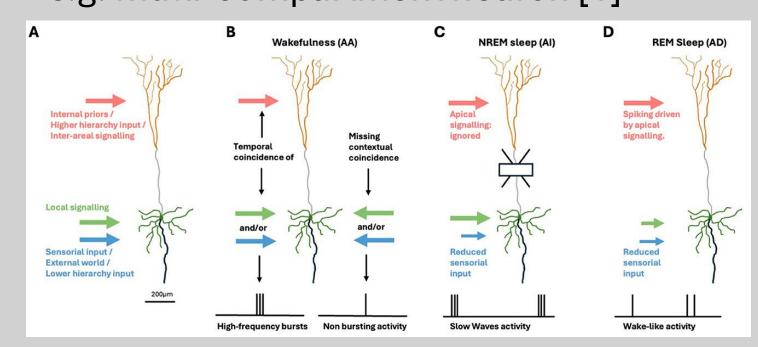
INFN Routing IP for low-latency inter-FPGA and intra-FPGA communication (< 1 us for up to 1 kB packets)





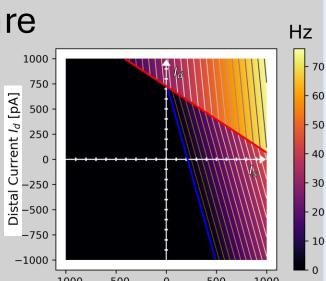


Adaptable neuron dynamics module to explore usage of biologically realistic neuron models e.g. multi-compartment neuron [1]



- Specialized in incremental learning.
- ThetaPlanes: piece-wise linear approximation of its transfer function for hardware implementation

 $\nu_F(I_s, I_d; \nu) = \Theta_\rho(1 - \Theta_H) \cdot \nu_- + \Theta_H \cdot \nu_+$ [1] Pastorelli et al. )2025) doi: 10.3389/fncom.2025.1566196

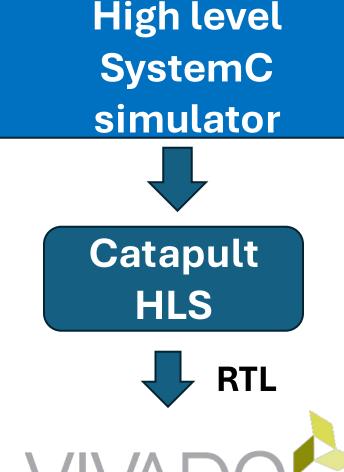


prototype with multiple levels of parallelism Configurable I/O cores for diverse input encoding

Modular, scalable, reconfigurable architecture

 Authomatized configurable synthesis workflow. Config » number of nodes, router ports (Aigor-I/O Core), neuron model, bit-field sizes, sim parameters, etc.

### Flexible development Workflow Based on High Level Synthesis



pipelines

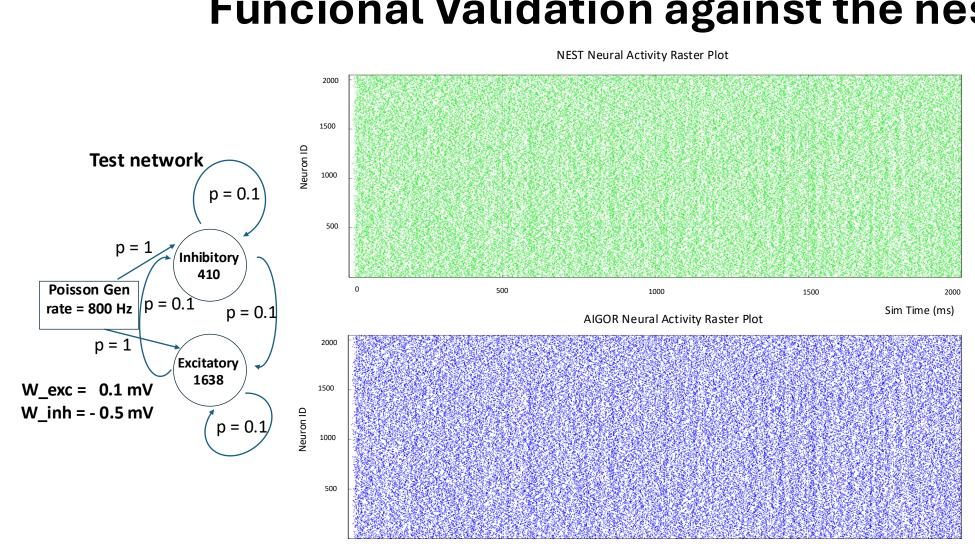




**FPGA** 

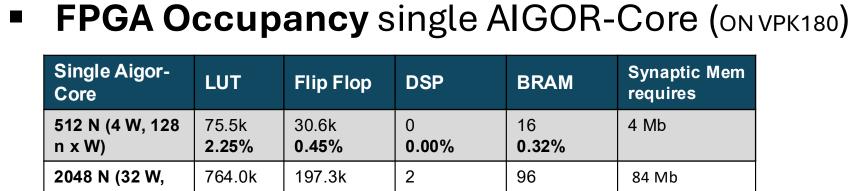
- SystemC (C++ based event-driven) architecture simulator
  - Validate and test architecture functionality
  - Modelling multi-FPGAs communication delays
  - SystemC allows higher or lower levels of abstraction (TLM, RTL-like)
- Catapult High Level Synthesis (HLS)
  - Direct transition from high-level simulator to RTL
  - Benefits: high level language, software emulation mode, early estimates on latency/throughput and FPGA-resources consumption  $\rightarrow$  relatively easy validation, reprogrammability, and debug

# Funcional Validation against the nest software simulator

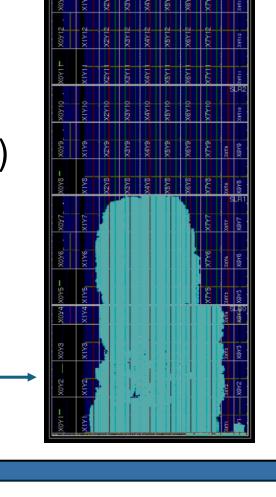


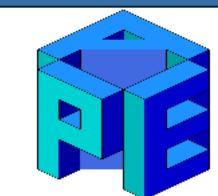
C++ simulator, with a Python interface, implementing a large number of models of biological neurons and synapses

nest:



0.01% 1.94% 2.93% N: tot. neurons, W: workers 708.5k 2048 N (16 W, 106.9k 84 Mb 0.24% 21.1% 1.59% 1.30% VPK-180 Total PL mem 994 Mb











n x W: neurons per worker

