# A System-on-Chip Based Hybrid Neuromorphic Compute (HNC) Node Architecture for Reproducible Hyper-Real-Time Simulations of Spiking Neural Networks

**Guido Trensch[1,3] and Abigail Morrison[1,2,3]**

[1]Simulation and Data Laboratory Neuroscience, [2]Institute of Neuroscience and Medicine (INM-6), JARA-Institute Brain Structure-Function Relationship (JBI-1 / INM-10), Research Centre Jülich, [3]Department of Computer Science 3 - Software Engineering, RWTH Aachen University

**JÜLICH** Forschungszentrum

ADVANCED COMPUTING ARCHITECTURES

## Goal

Benefiting from the continued advances in semiconductor technology, in recent years, programmable device technology and tools have greatly increased.

**Proof of Concept:** Prototypical implementation of an AMD Xilinx System-on-Chip (SoC) based hybrid software-hardware architecture approach for a neuromorphic compute node capable of meeting the high demands for modeling and simulation in neuroscience.

## HNC Node High-Level Architecture
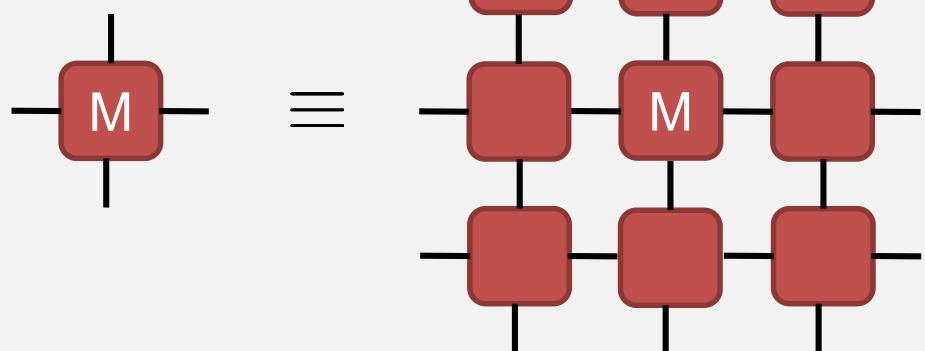


**Simulation paradigm**
- Hybrid strategy for time-discrete neural network simulations of point neuron models:
  - time-driven neuron state update (blue data paths) at fixed intervals, $\Delta t = 0.1ms$
  - event-driven synapse update (red data paths)

**Performance**
- Exploiting the tight coupling of an Application Processing Unit (APU) with a Field Programmable Gate Array (FPGA) located on the same chip.
- Off-loading of performance critical algorithms to programmable logic.
- Parallelization by distributing the computational load over 16 processing units (P1, P2, …, P16).
- Data locality of state variables by storing them in fast on-chip block RAM memories (BRAMs).
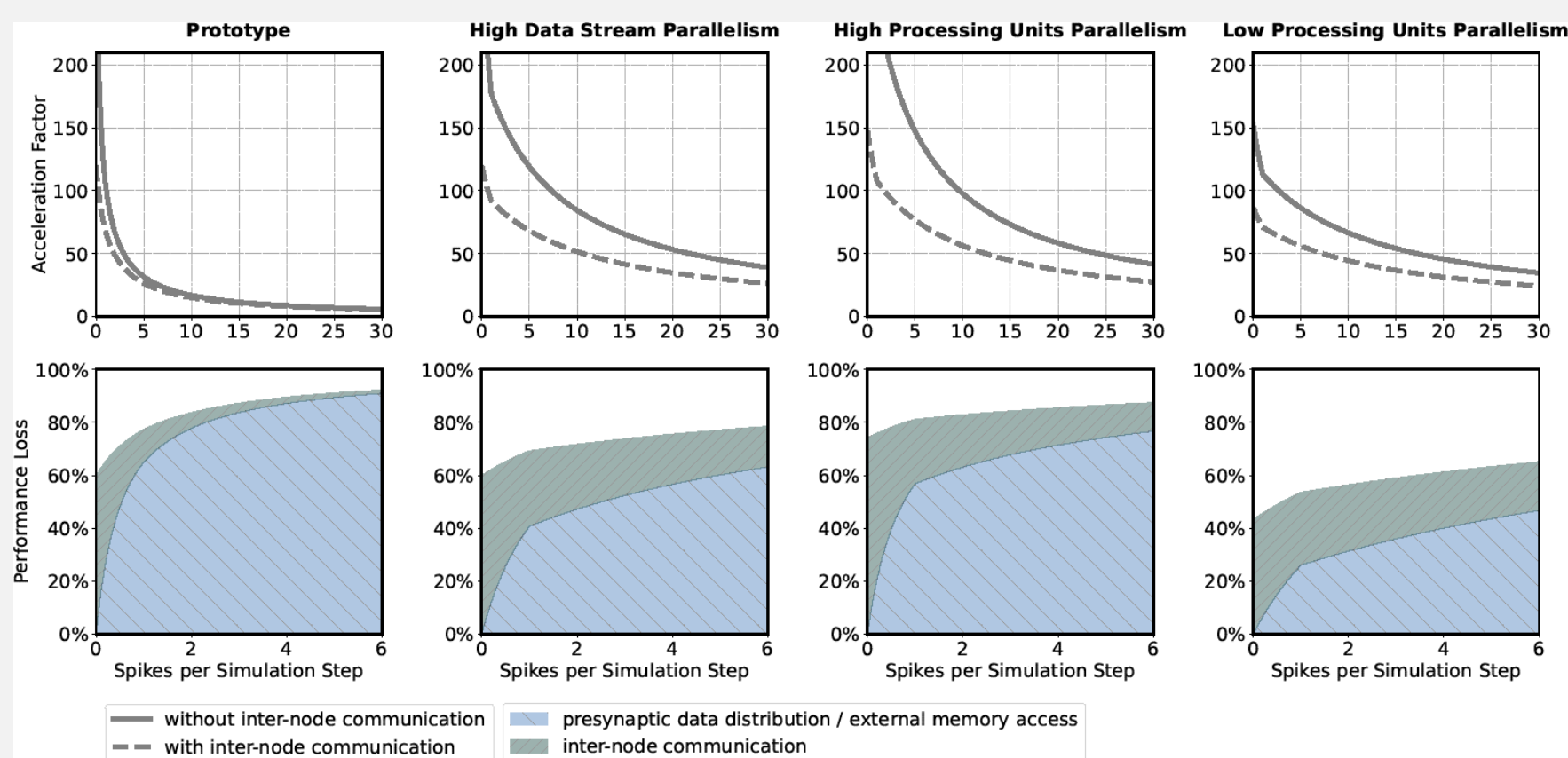- Latency hiding by exploiting the true-dual port capability of BRAMs.

**Flexibility**
- ODE pipeline module can be replaced to implement different neuron and synapse models.
- Flexibility in the choice of data types.
- Connectivity data is stored in external memory, thus synaptic weights are adjustable and accessible to the APU.
- Non-critical tasks are executed in software by the APU.

## Workload Model

**From the workload perspective, it can be considered equivalent.**

**small** network **high** average firing rate

**large** network **low** average firing rate

$$\bar{\nu} = \frac{1}{N} \sum_N \frac{n_{sp}(T)}{T}$$

**Indicator for the computational workload: average number of spikes per time step $k$**

$$\bar{\nu}_k = \frac{h}{T} \sum_N n_{sp}(T) = N\bar{\nu}h$$

| Parameter | | HNC Node |
|---|---|---|
| number of compute nodes | $M$ | N/A |
| total number of neurons | $N$ | N/A |
| number of neurons per compute node | $N^M$ | max. 1024 |
| connection probability | $\epsilon$ | < 0.1 |
| presynaptic neuron's max. number of connections per node | $C^M = \epsilon N^M = \dfrac{\epsilon N}{M}$ | max.128 |
| spike count of neuron $n$ in interval $T$ | $n_{sp}(T)$ | N/A |
| temporal resolution of the simulation | $h$ | 0.1 ms |

## Performance Model

**Acceleration factors derived from operating latencies.**

Max. acceleration factor (single node): $\quad F_S^{MAX} = \dfrac{h f_{clk}}{L_\Sigma}$

Without communication (single node):
$$F_S(\bar{\nu}_k) = \begin{cases} \dfrac{h f_{clk}}{\bar{\nu}_k(L_\Sigma^{SE} + L_{DS}) + (1 - \bar{\nu}_k)L_\Sigma} & \text{if } \bar{\nu}_k < 1 \\ \dfrac{h f_{clk}}{L_\Sigma^{SE} + \bar{\nu}_k L_{DS}} & \text{otherwise} \end{cases}$$

With communication (cluster):
$$F_C(\bar{\nu}_k) = \begin{cases} \dfrac{h f_{clk}}{\bar{\nu}_k(L_\Sigma^{SE} + L_{DS} + \alpha L_{COM}) + (1 - \bar{\nu}_k)L_\Sigma + L_{COM}} & \text{if } \bar{\nu}_k < 1 \\ \dfrac{h f_{clk}}{L_\Sigma^{SE} + \bar{\nu}_k(L_{DS} + \alpha L_{COM}) + L_{COM}} & \text{otherwise} \end{cases}$$

## System-Level Architecture



**DMA** Direct Memory Access
**RB** Ring Buffer
**SVB** State Variables Buffer
**PRNG** Pseudo Random Number Generator
**GIC** Global Interrupt Controller
**OCM** On-Chip Memory
**P1..16** Processing Units

XCZ7045

### Software system executed on the APU



- Orchestrates the overall node operation.
- A minimal C-API provides `Create()`, `Connect()`, and `Simulate()` function calls.
- Implemented as bare-metal application in C.

## Single Node Performance

**1000 neurons two-population Izhikevich network.**



$f_{clk} = 200MHz$

**HNC Node**

Data locality of state variables

External memory access (~1.8 GB/s)

Simulate workload: consecutive simulation runs of 5 min simulated biological time with an increasing external offset current. $\quad i_{ext} = \{-3.0pA, .., +100pA\}$

## Performance Characteristics



Prototype | High Data Stream Parallelism | High Processing Units Parallelism | Low Processing Units Parallelism

| Parameters | Prototype | High Data Stream Parallelism | High Proc. Units Parallelism | Low Proc. Units Parallelism |
|---|---|---|---|---|
| number of data streams, $DS$ | 2 | 16 | 16 | 16 |
| data stream latency, $L_{DS}$ | 110 | 14 | 14 | 14 |
| number processing units, $P$ | 16 | 16 | 32 | 8 |
| number of neurons per processing unit, $N^P$ | 64 | 64 | 32 | 128 |
| ODE pipeline iteration latency, $IL_N$ | 64 | 64 | 32 | 128 |
| **Acceleration factors w/o communication** | | | | |
| maximum, $F_S^{MAX} = F_S(\nu_k = 0)$ | 298.5 | 298.5 | 571.4 | 152.7 |
| low workload, $F_S(1.0)$ | 104.7 | 177.0 | 246.9 | 113.0 |
| medium workload, $F_S(10.0)$ | 16.9 | 84.5 | 97.7 | 66.5 |
| high workload, $F_S(20.0)$ | 8.8 | 52.4 | 58.4 | 45.6 |
| **Acceleration factors with communication** | | | | |
| maximum, $F_C^{MAX} = F_C(0)$ | 119.8 | 119.8 | 148.1 | 86.6 |
| low workload, $F_C(1.0)$ | 67.6 | 91.7 | 107.5 | 70.9 |
| medium workload, $F_C(10.0)$ | 15.0 | 51.7 | 56.4 | 44.4 |
| high workload, $F_C(20.0)$ | 8.1 | 34.8 | 36.9 | 31.3 |

without inter-node communication
with inter-node communication
presynaptic data distribution / external memory access
inter-node communication

## Prototyping Platform

**AMD Xilinx Zynq®-7000 SoC ZC706 Development Board**



Dual-core ARM Cortex-A9 CPU (up to 1GHz) + FPGA (250K logic cells, 19.1Mb Block RAM, 900 DSP slices)

1GB DDR3 external memory

Contact: Guido Trensch, g.trensch@fz-juelich.de