

Hierarchical Network-on-Chip and Traffic Compression for Spiking Neural Network Implementations

Snaider Carrillo, Jim Harkin, Liam McDaid
University of Ulster, Magee Campus

Sandeep Pande, Seamus Cawley, Brian McGinley, Fearghal Morgan
National University of Ireland, Galway Campus

Outline

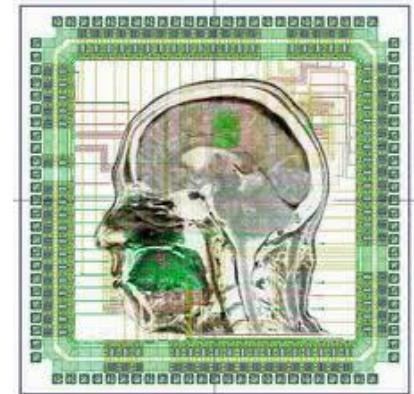
- Motivation and Challenges
- Hierarchical NoC EMBRACE Architecture
- Performance Analysis
- Take-home Message & Future Work



Motivation: Engineer & Neuroscientist

Neural processing systems.....Taking inspiration from the biology.....to deploy a new computer architecture paradigm !!!

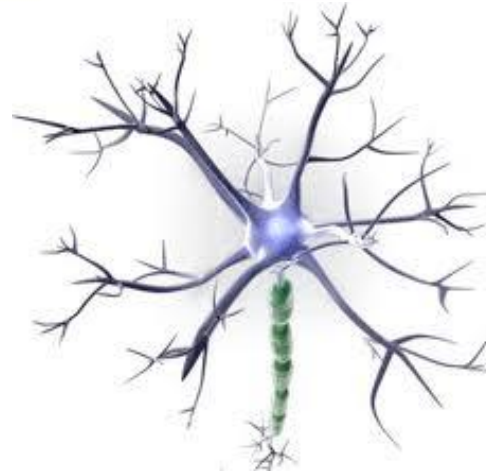
- **An Engineering point of view....**
 - Pattern recognition + Low power consumption
 - Fault-tolerant computers +Self repairing systems
- **A Neuroscientist point of view....**
 - Faster large-scale neural network simulations
 - Ultimately, to learn a bit more about how the human brain works



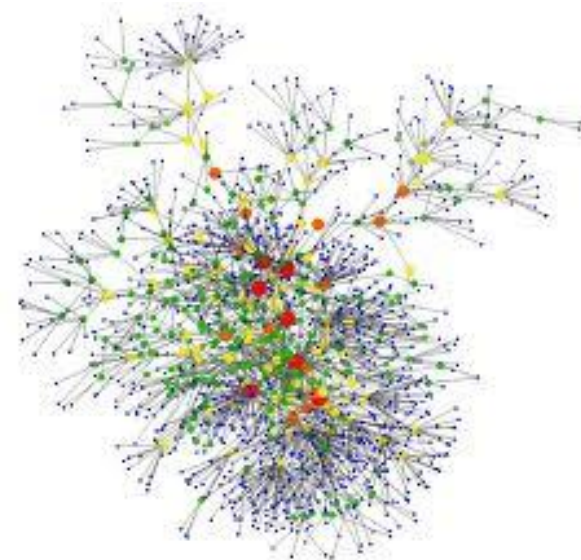
Neuron Interconnection: The big challenge

A human brain contains in average...

- 10^{11} neurons
- 10^{15} synapses
- 1:1000 Fan in/out connection ration



We need a highly optimised architecture to overtake this challenge....



Previous Work

Blue Brain Project [Markram'03]

- IBM BlueGene/L supercomputer

SpiNNaker [Furber'06]

- Embedded ARM processors + NoC interconnection

SYNAPSE Project [Modha'11]

- Digital neurons + Crossbar fabric

Neurogrid [Boahen'09]

- Analogue neurons + on-chip routers

FACETS [Schemmel'05]

- Analogue neurons + hierarchical intra-wafers buses

....However, there is still room for improvement 😊

Key Research Problem

*...How to **interconnect** a large number of **spiking neurons** in a **network** fashion **efficiently**?...*

Efficiently?... a trade-off between

- *Scalability*
- *Area utilisation*
- *Power consumption*
- *Throughput*
- *Synapse/neuron ratio*

.....And what about hardware acceleration !!



Outline

- Motivation and Challenges
- Hierarchical NoC EMBRACE Architecture
- Performance Analysis
- Take-home Message & Future Work



EMulating BBiologically-inspiRed Architectures in hardwarE (EMBRACE)

Self-repairing Embedded
Information Processing Systems

Accelerated Exploration Platform
for Neuro-degenerative Diseases

EMBRACE

Electronic
Biological Cells

- CMOS Synapse
- Neuron cell

Low-level

Electronic
Storage

- Weight storage
- re-programming architectures

Interconnect

- NoCs
- Adaptive routers
- Fault detection

Computational
Models

- Astrocyte models
- Self repair models
- Learning models

Tools

- Network Builder
- Programming
- Analysis tool

High-level

- Ulster
- University of Liverpool (S Hall)

- Ulster
- NUI Galway (F Morgan)

- Ulster
- University of Cardiff
(Prof. V Cruneli)

- NUI Galway (F Morgan)

EMulating BBiologically-inspiRed Architectures in hardwarE (EMBRACE)

Self-repairing Embedded
Information Processing Systems

Accelerated Exploration Platform
for Neuro-degenerative Diseases

EMBRACE

**Electronic
Biological Cells**

- CMOS Synapse
- Neuron cell

Low-level

**Electronic
Storage**

- Weight storage
- re-programming architectures

Interconnect

- NoCs
- Adaptive routers
- Fault detection

**Computational
Models**

- Astrocyte models
- Self repair models
- Learning models

Tools

- Network Builder
- Programming
- Analysis tool

High-level

- Ulster
- University of Liverpool (S Hall)

- Ulster
- NUI Galway (F Morgan)

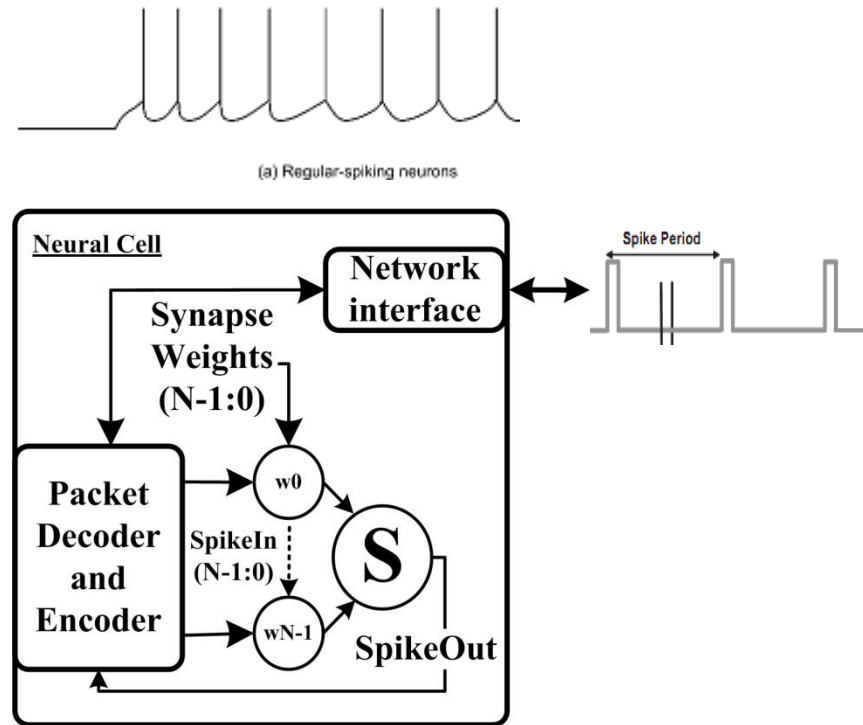
- Ulster
- University of Cardiff
(Prof. V Cruneli)

- NUI Galway (F Morgan)

EMBRACE Neural Cell

■ Provides:

- An analogue **point neuron** (Leaky Integrate & Fire model)
- Its correspondent **synapse cells** (Dynamic Synapses)
- A **packet decoder/encoder**
- A **network interface** to communicate with **digital NoC router**



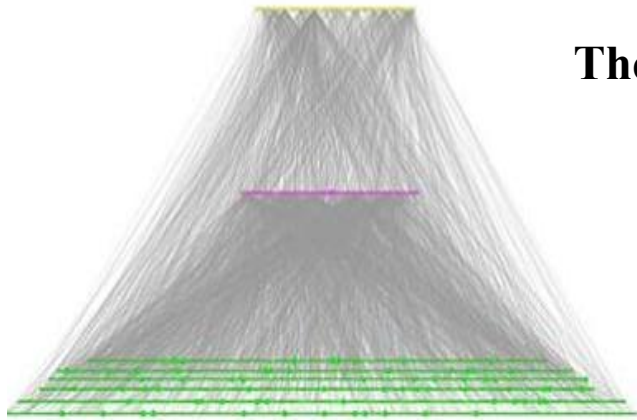
On-going EPSRC project between:

- University of Ulster
- University of Liverpool (S Hall)

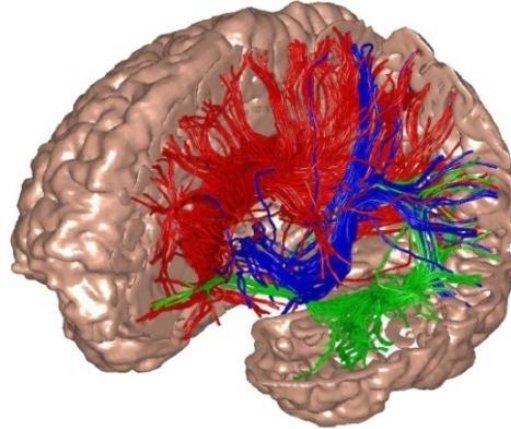
L. McDaid, S. Hall, and P. Kelly, "A programmable facilitating synapse device," in *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, 2008, pp. 1615-1620.

Hierarchical Topology: Taking inspiration from the biology.....

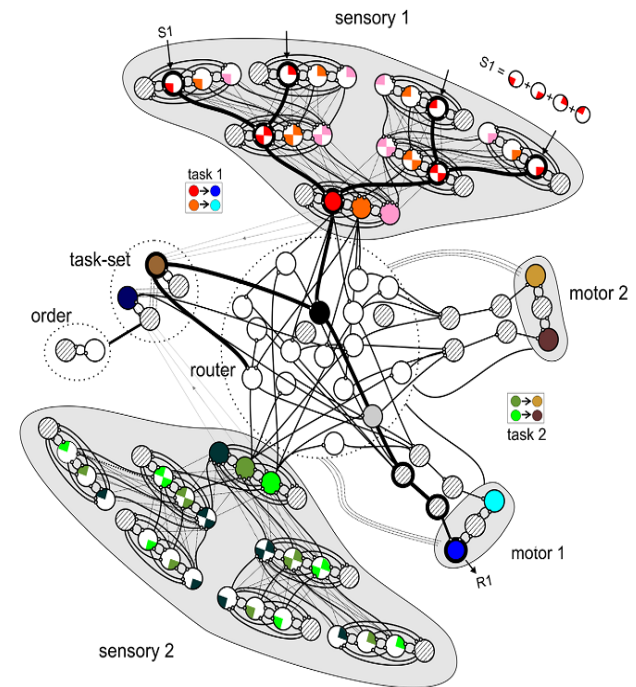
E. coli transcriptional regulatory network



The hierarchical topology of
the *E. Coli* (Yan et al. 2010)



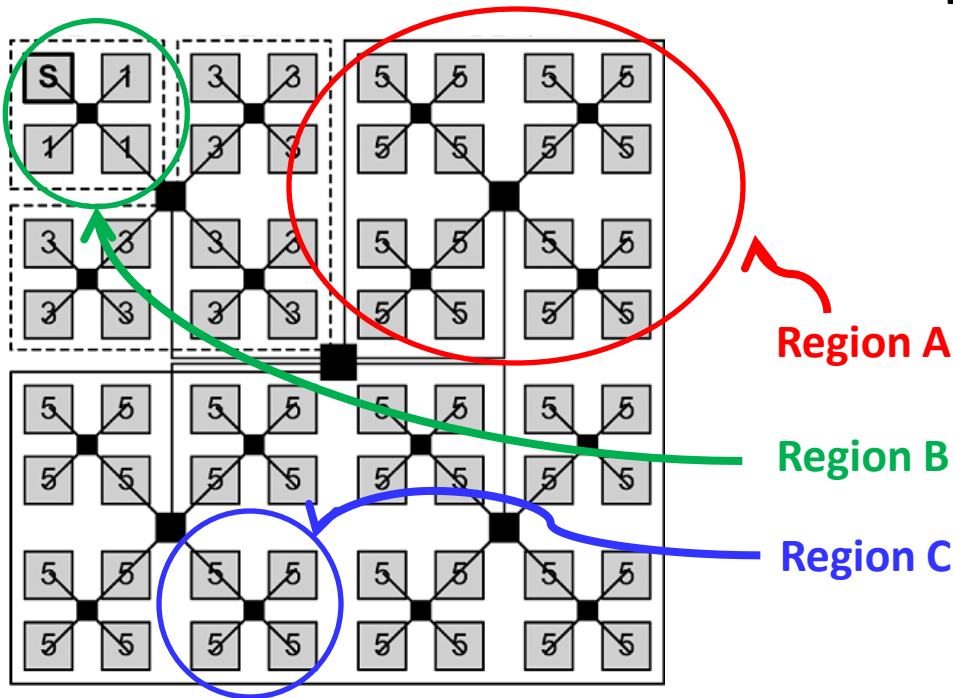
The brain is a 3D structure !!



A Schematic representation of a cluster of
neurons (Zylberberg et al. 2010)

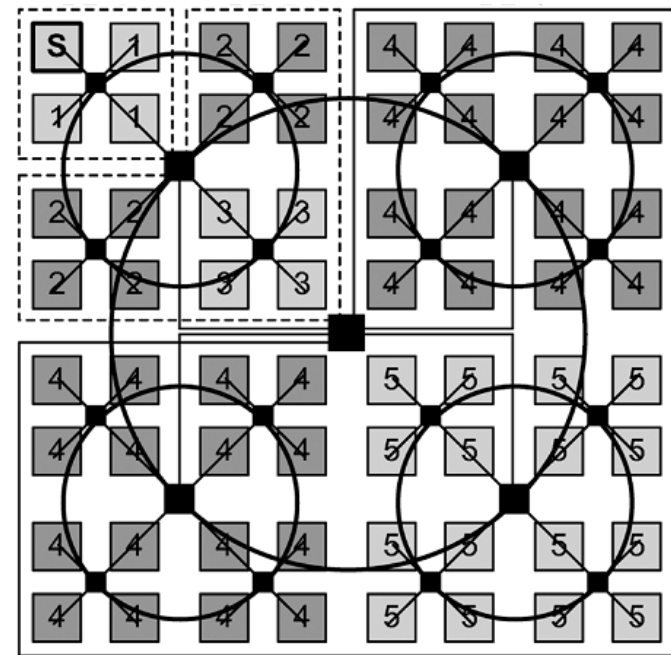
Hierarchical Topology: and also from the NoC community !!

- Hierarchical NoCs (**H-NoCs**) exploit the concept of **region-based routing**.



Hierarchical star [1]

- **Virtual regions** or facilities are used to allocate resources that process either **local or global traffic**.

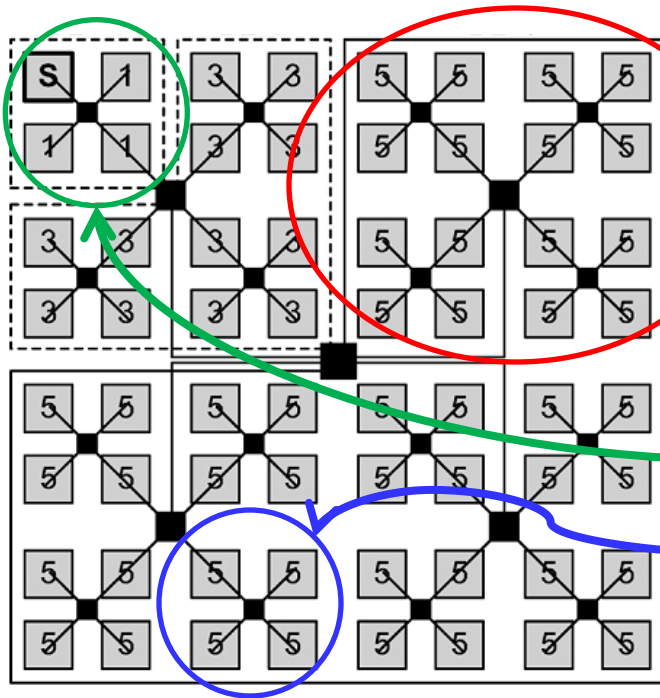


Hierarchical star + ring [1]

[1] J.-Y. Kim, J. Park, S. Lee, M. Kim, J. Oh, and H.-J. Yoo, "A 118.4 GB/s Multi-Casting Network-on-Chip With Hierarchical Star-Ring Combined Topology for Real-Time Object Recognition," *IEEE Journal of Solid-State Circuits*, vol. 45, no. 7, pp. 1399-1409, Jul. 2010

Hierarchical Topology: and also from the NoC community !!

- Hierarchical NoCs (**H-NoCs**) exploit the concept of **region-based routing**.

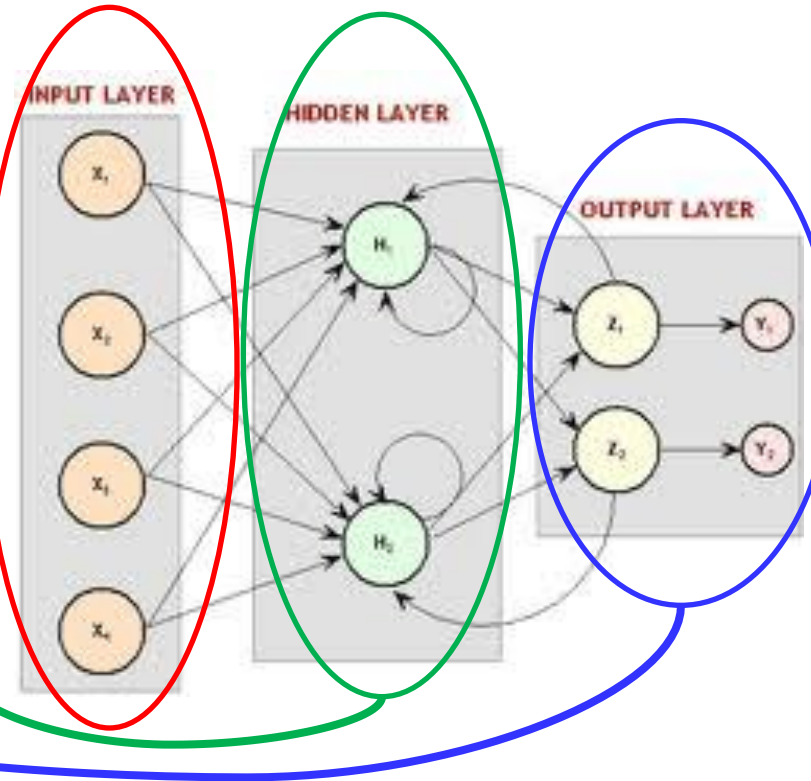


Hierarchical star [1]

Region A

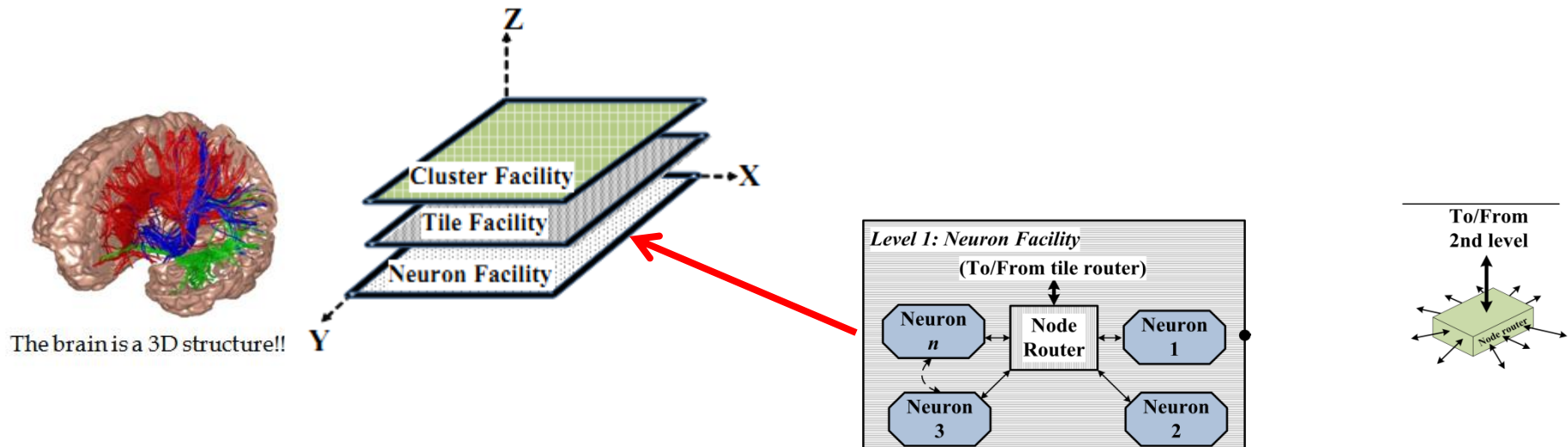
Region B

Region C

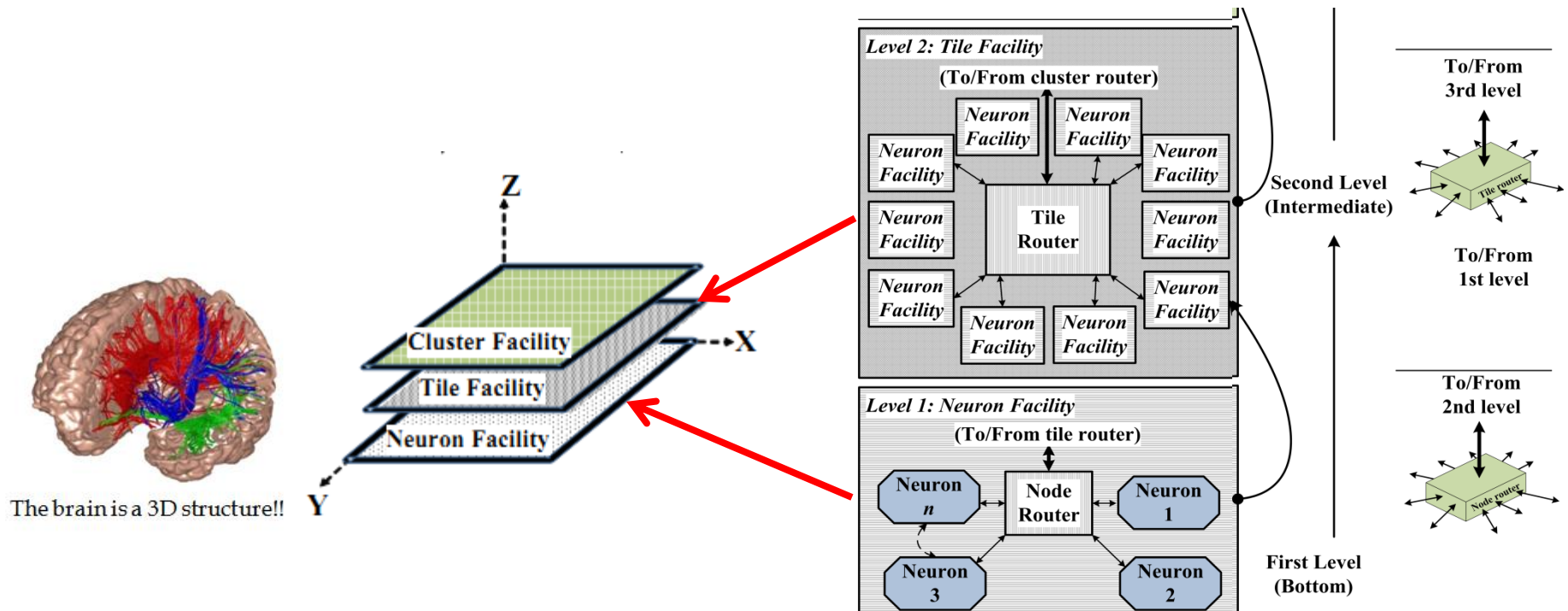


[1] J.-Y. Kim, J. Park, S. Lee, M. Kim, J. Oh, and H.-J. Yoo, "A 118.4 GB/s Multi-Casting Network-on-Chip With Hierarchical Star-Ring Combined Topology for Real-Time Object Recognition," *IEEE Journal of Solid-State Circuits*, vol. 45, no. 7, pp. 1399-1409, Jul. 2010

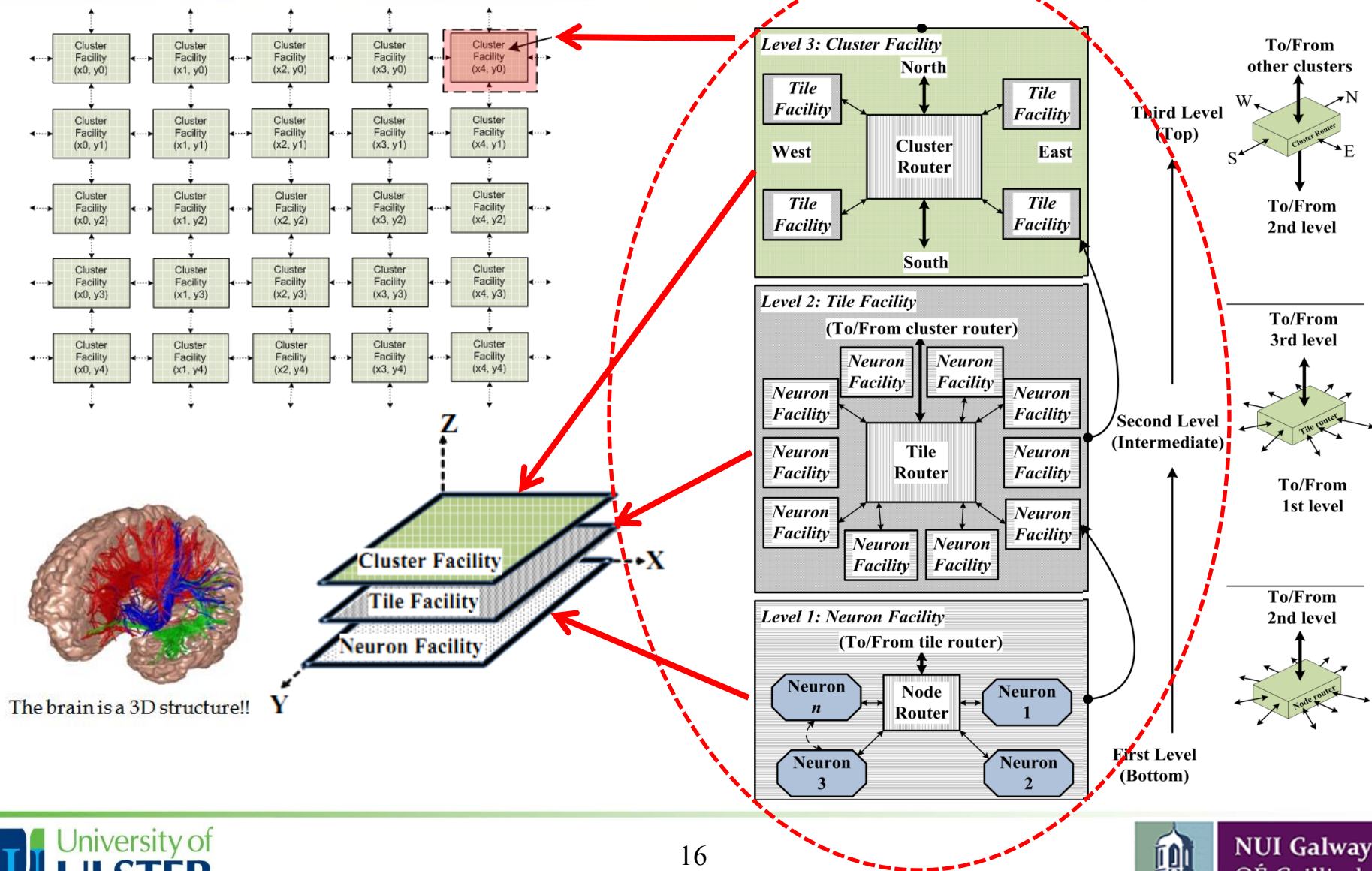
EMulating BBiologically-inspiRed Architectures in hardware (EMBRACE): H-NoC Approach



EMulating BBiologically-inspiRed Architectures in hardware (EMBRACE): H-NoC Approach



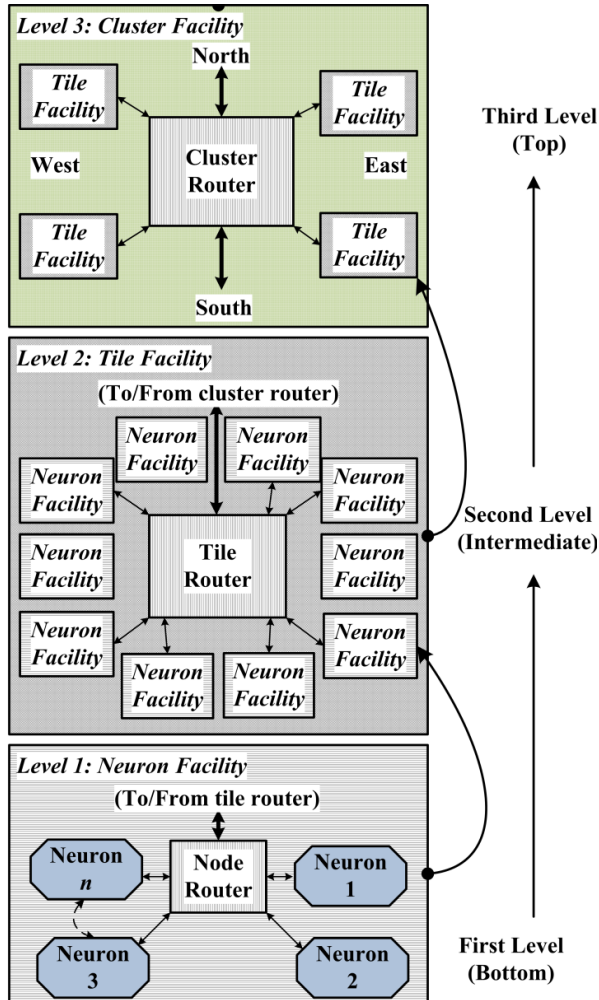
EMulating Biologically-inspired Architectures in hardware (EMBRACE): H-NoC Approach



EMBRACE: H-NoC Architecture

One Cluster Facility contains:

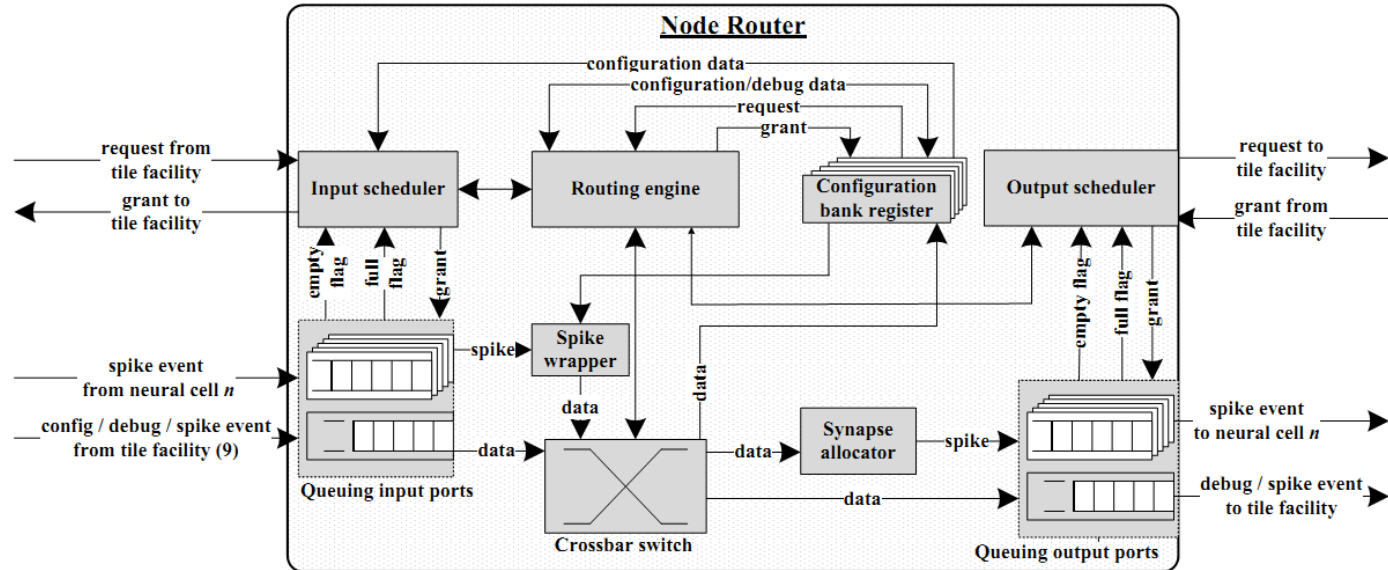
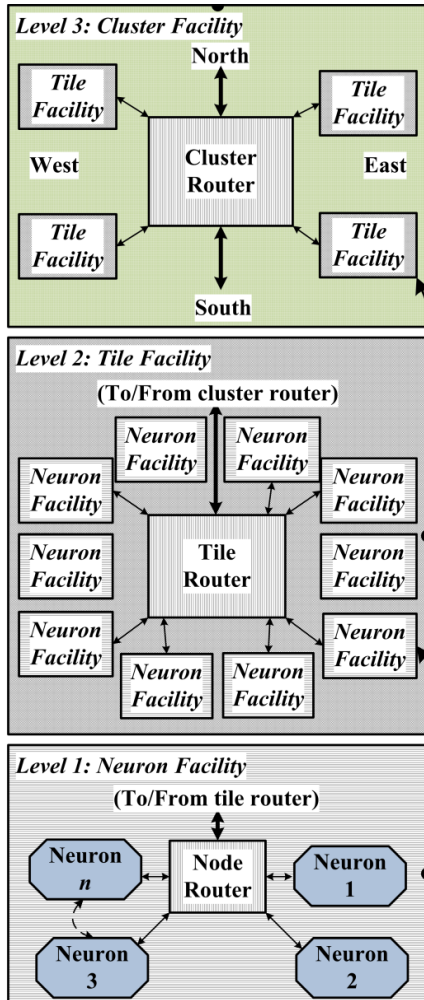
+ 1 Cluster NoC router
4 Tile NoC routers
40 Node NoC Router



A total of **45 NoC Router**
to interconnect **400 neurons**
....This is just an initial density 😊

Carrillo, S., et al., "Advancing Interconnect Density for Spiking Neural Network Hardware Implementations using Traffic-Aware Adaptive Network-on-Chip Routers". Neural Networks, Vol 33, pp. 42-57, September 2012.

Neuron Facility – @Bottom-Level



Layout for the configuration bank register

Neuron facility address register

Address range [3:0]

Spike event timer register

Scaling factor range [13:0]

Communication protocol register

IT_CM/NB_MC/NB_BC/DB [3:0]

Spike event reception register

Neurons position within node [9:0]

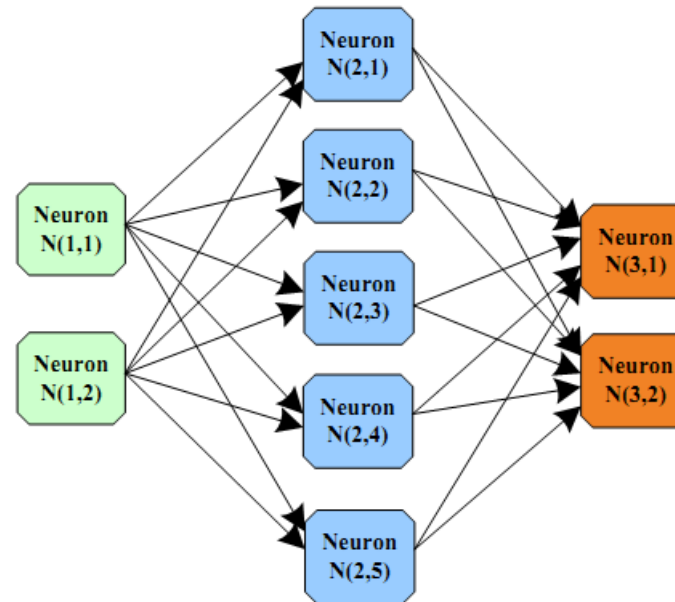
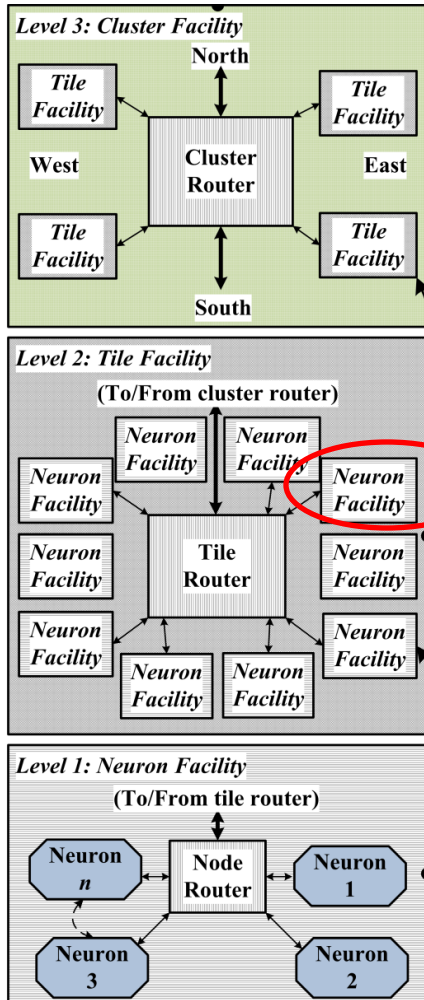
Target neuron population register

Target X-cluster addr [29:22]	Target Y-cluster addr [21:14]
Target tile addr [13:10]	Target node addr [9:0]

NB_MC: neighbour multicast DB: debug

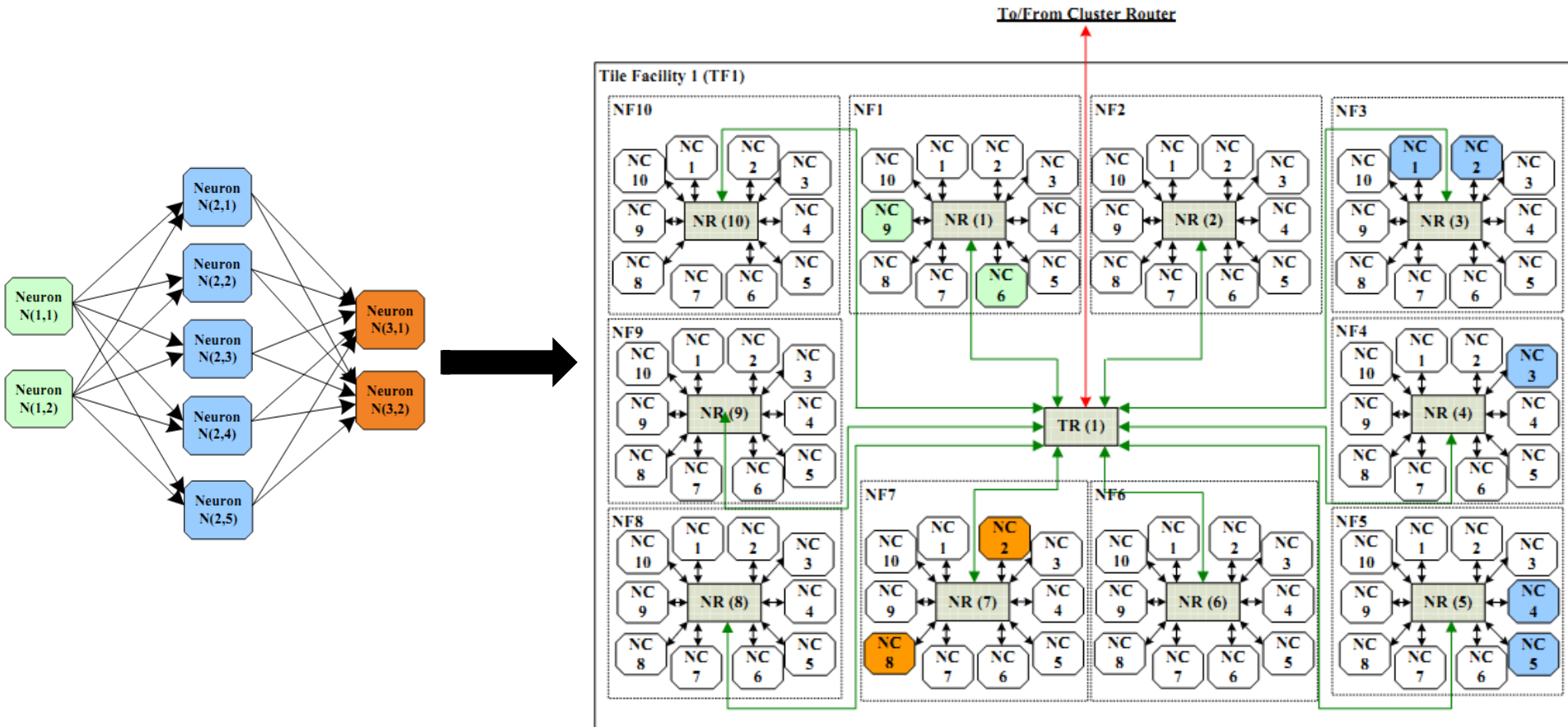
NB_BC: neighbour broadcast IT_CM: internal communication

H-NoC Architecture: Example Scenario

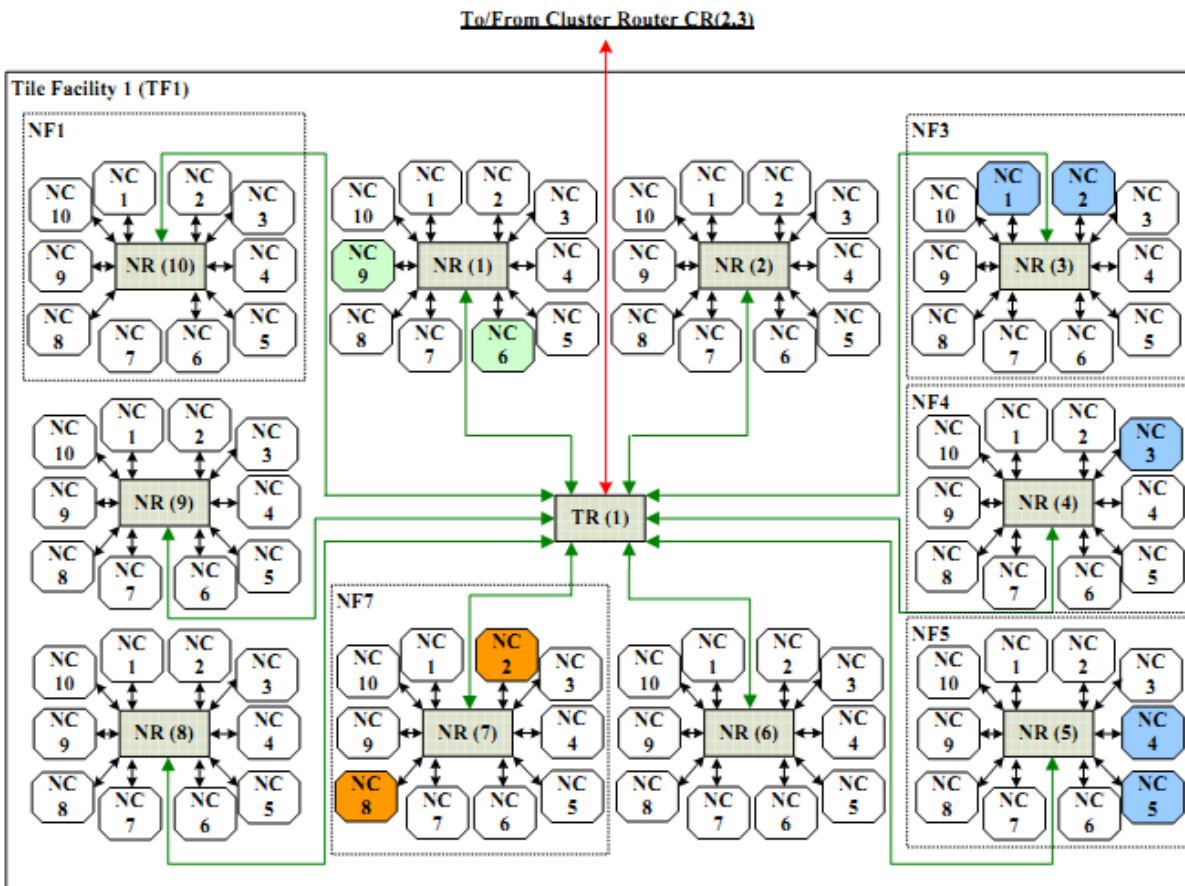


2x5x3 Feed-forward neural network

H-NoC Architecture: Example Scenario



H-NoC Architecture: Example Scenario



Node Router NR(1)

<u>Address Register</u>	01h	<u>Comm. Protocol</u>	0Fh
<u>Reception Register</u>	000h	<u>Timer Register</u>	0001h
<u>Target Register</u>	02h	03h	1h 01Ch

Node Router NR(3)

<u>Address Register</u>	03h	<u>Comm. Protocol</u>	0Fh
<u>Reception Register</u>	003h	<u>Timer Register</u>	0001h
<u>Target Register</u>	02h	03h	1h 040h

Node Router NR(4)

<u>Address Register</u>	04h	<u>Comm. Protocol</u>	0Fh
<u>Reception Register</u>	004h	<u>Timer Register</u>	0001h
<u>Target Register</u>	02h	03h	1h 040h

Node Router NR(5)

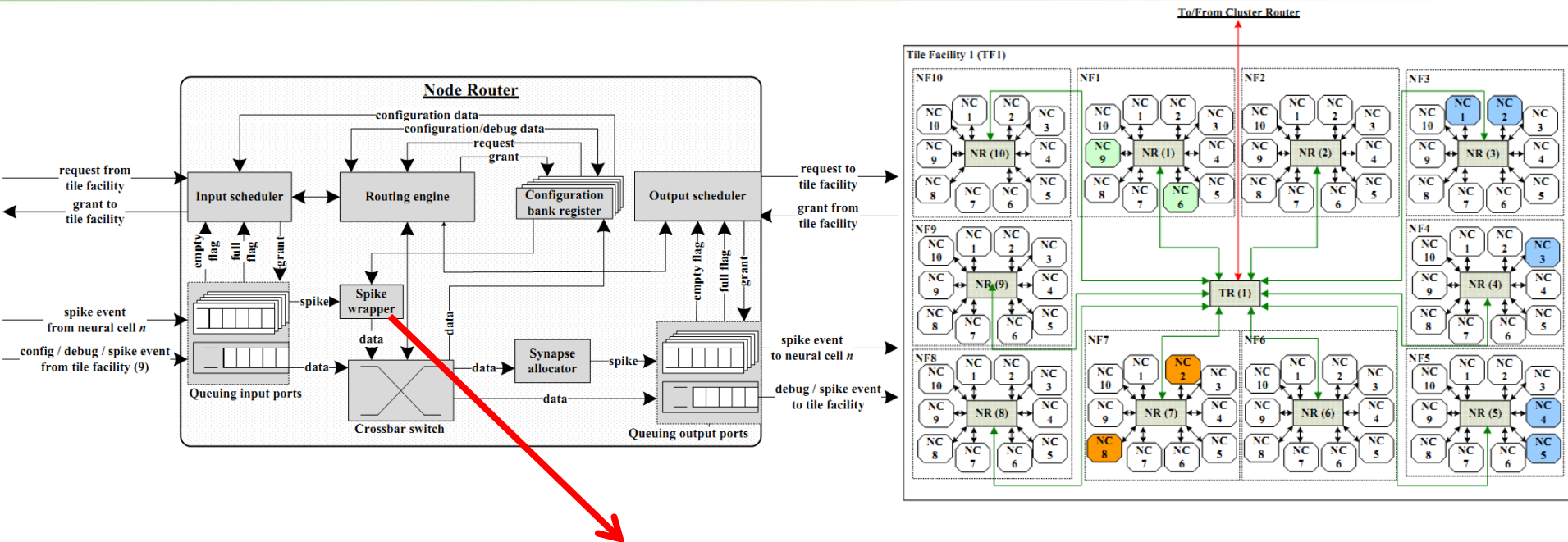
<u>Address Register</u>	05h	<u>Comm. Protocol</u>	0Fh
<u>Reception Register</u>	018h	<u>Timer Register</u>	0001h
<u>Target Register</u>	02h	03h	1h 040h

Node Router NR(7)

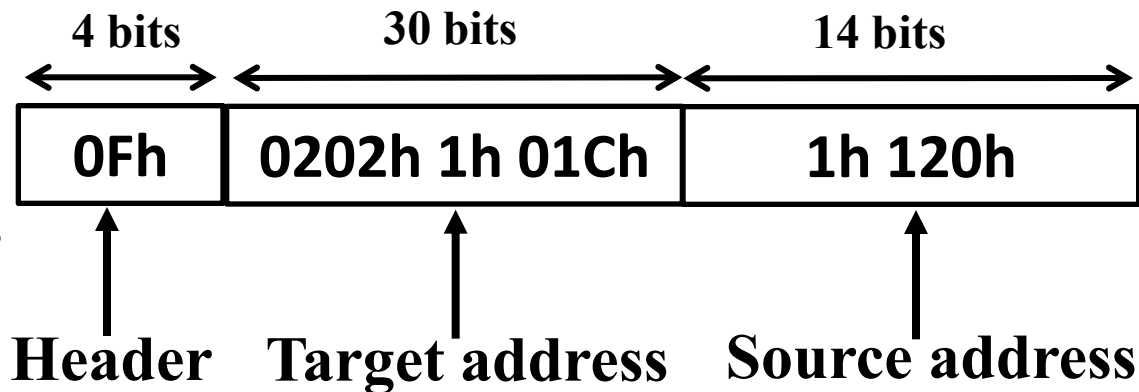
<u>Address Register</u>	07h	<u>Comm. Protocol</u>	0Fh
<u>Reception Register</u>	082h	<u>Timer Register</u>	0001h
<u>Target Register</u>	00h	00h	0h 000h

H-NoC Architecture: Example Scenario

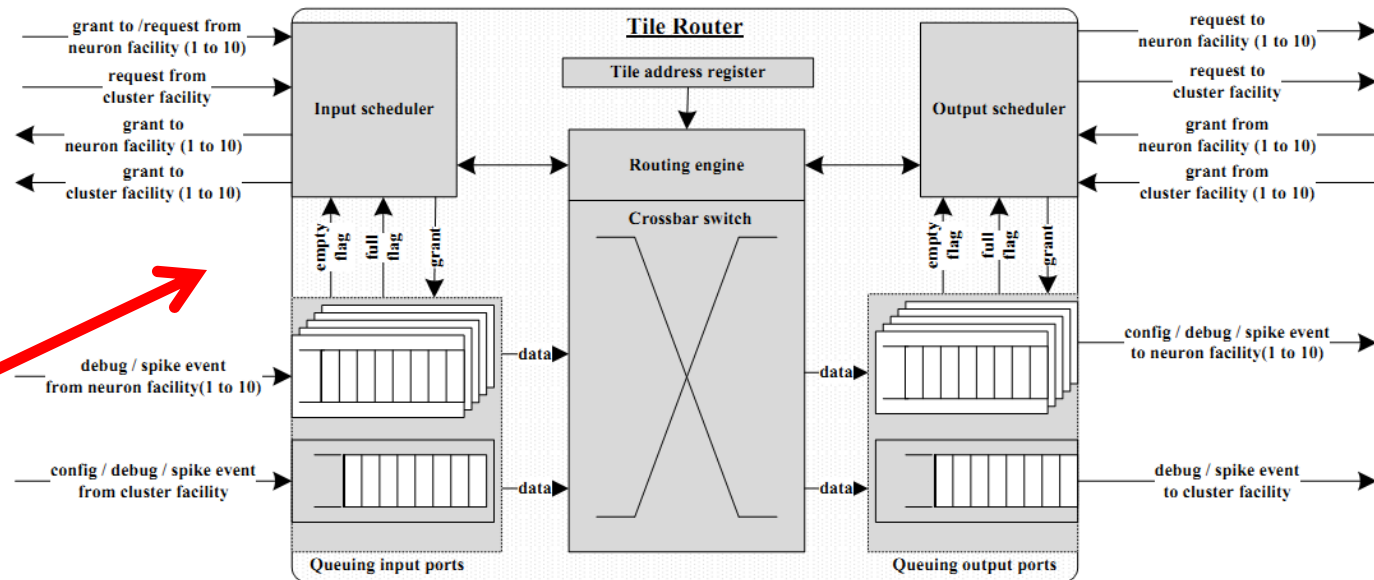
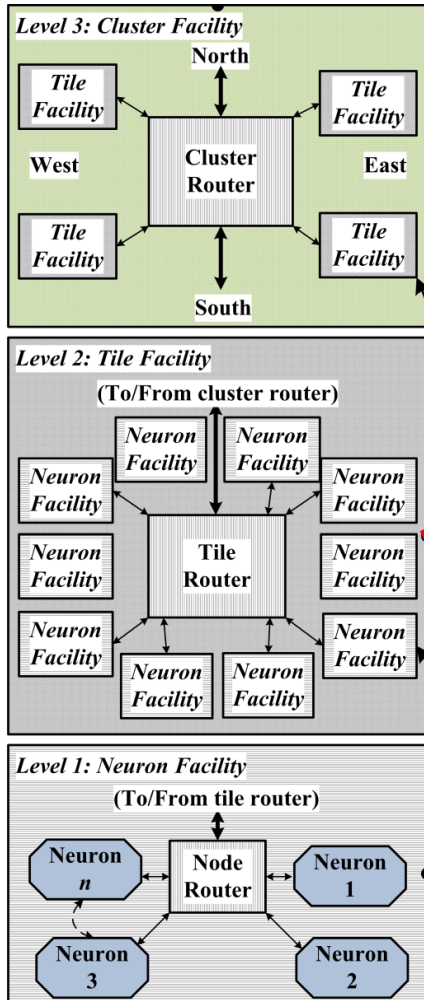
On-chip Comm: Spike event generation



Packet generated when
Input neurons #6 and # 9
are generating spike events



Tile Facility – @Mid-Level

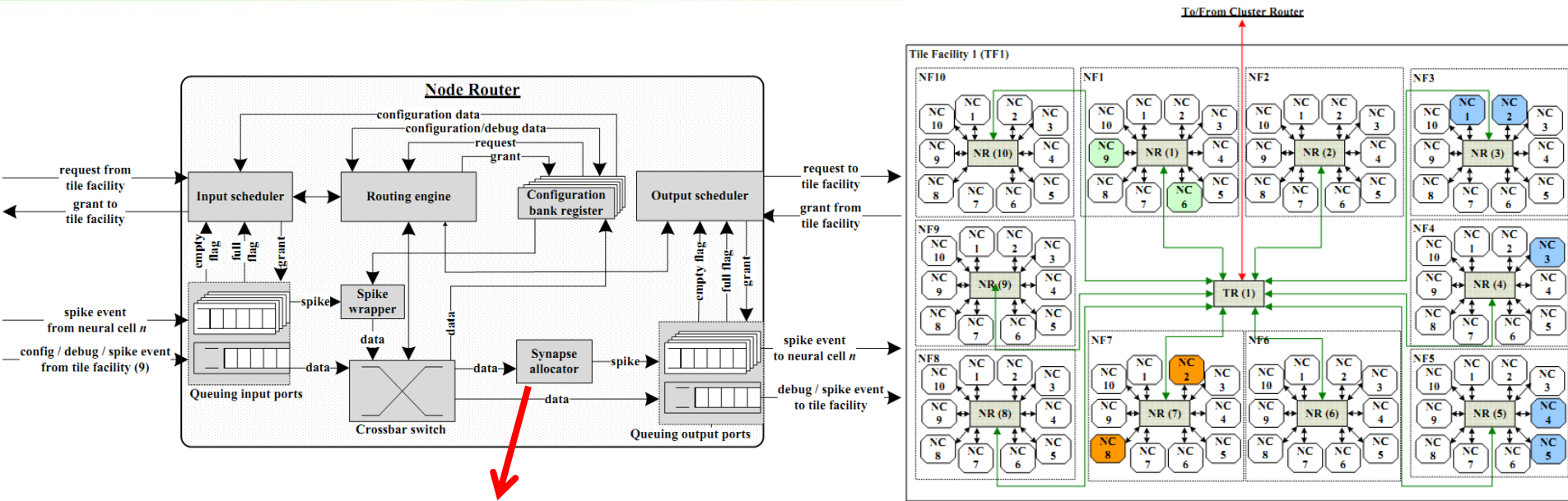


It's the arbitration point for NoC packets coming from the Bottom & Top Levels !!

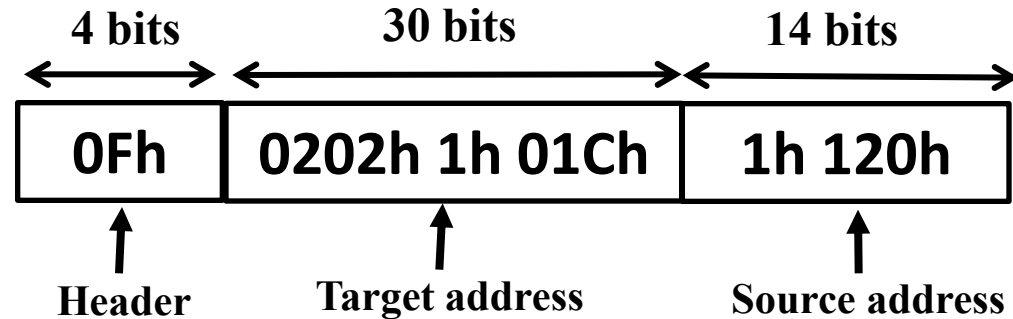
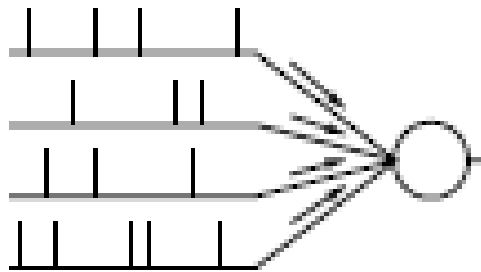
Distributed parallel datapath to handle multiple incoming spike events !!

H-NoC Architecture: Example Scenario

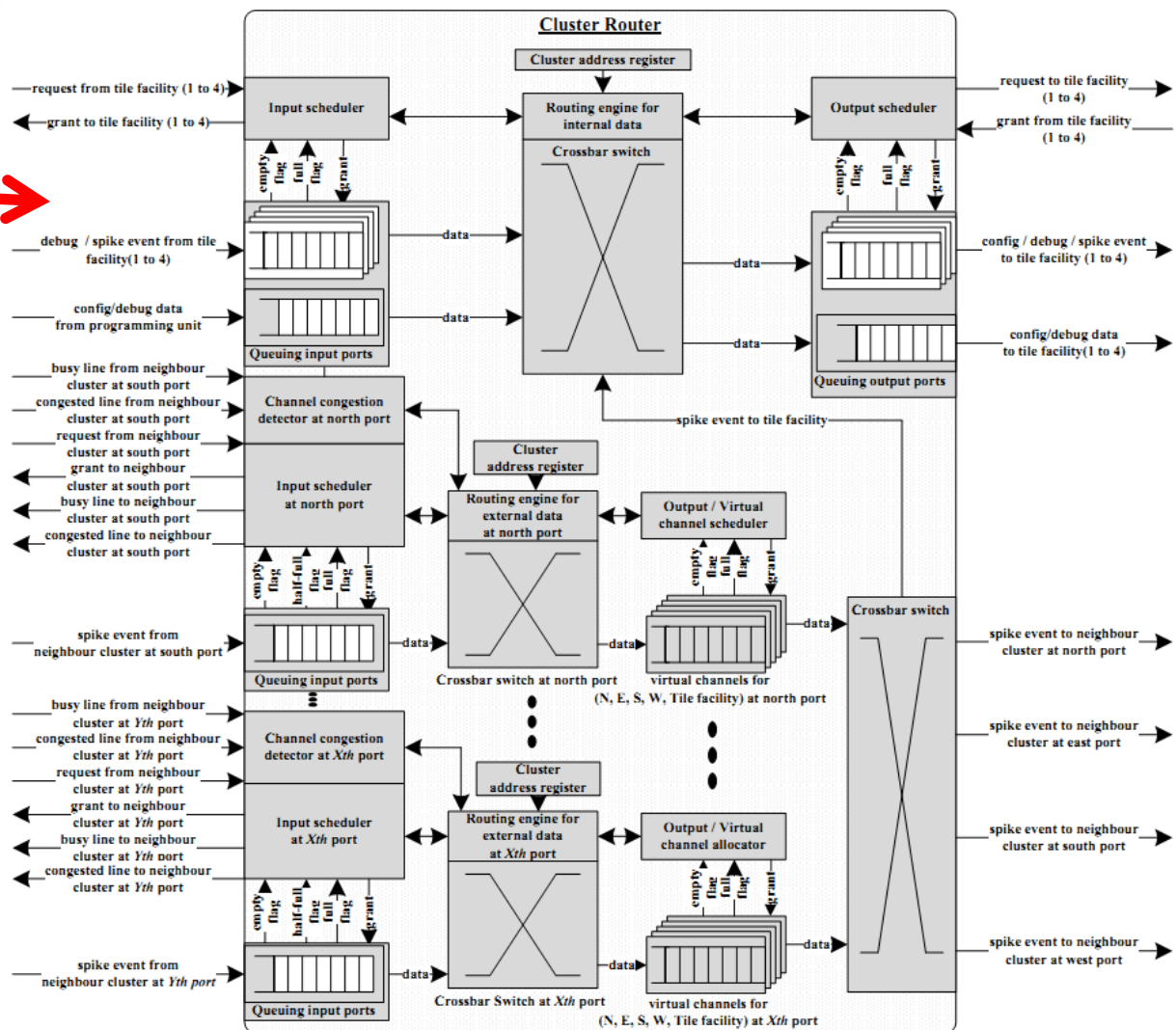
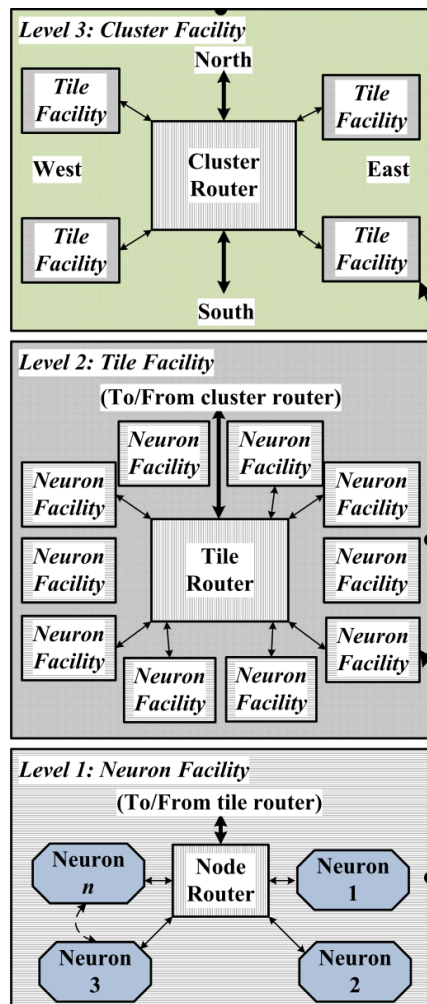
On-chip Comm: Spike event absorption



$$As = \eta - [A_{neuron} + (A_{neuron_facility} \times \eta)]$$



Cluster Facility – @Top Level



On-chip Communication Protocols & Free Look-up Table Approach

Layout for the configuration bank register

Neuron facility address register

Address range [3:0]

Spike event timer register

Scaling factor range [13:0]

Communication protocol register

IT_CM/NB_MC/NB_BC/DB [3:0]

Spike event reception register

Neurons position within node [9:0]

Target neuron population register

Target X-cluster addr [29:22]

Target Y-cluster addr [21:14]

Target tile addr [13:10]

Target node addr [9:0]

NB_MC: neighbour multicast DB: debug

NB_BC: neighbour broadcast IT_CM: internal communication

Packet layout for the configuration mode

Header [47:44]	Target address* [43:30]	Configuration data [29:0]
----------------	-------------------------	---------------------------

Packet layout for the debug mode

Header [47:44]	Debug information [43:14]	Source address [13:0]
----------------	---------------------------	-----------------------

Packet layout for the execution mode

Header [47:44]	Target address [43:14]	Source address [13:0]
----------------	------------------------	-----------------------

Target address field

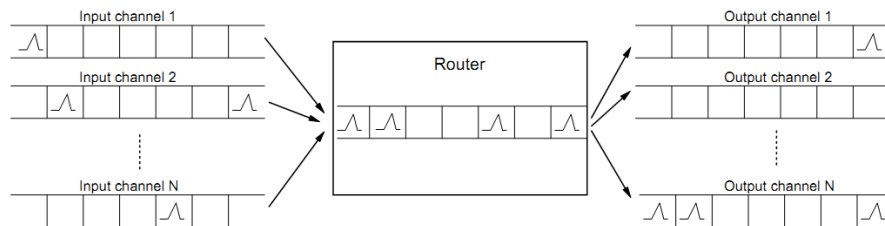
Target X-cluster addr [43:36]	Target Y-cluster addr [35:28]
Target tile addr [27:24]	Target node addr [23:14]

Source address field

Node address [13:10]	Neurons position within node [9:0]
----------------------	------------------------------------

Header Information

Header Information	Value
Cluster address register	0x00h
Tile address register	0x01h
Node address register	0x02h
Spike event timer	0x03h
Communication protocol register	0x04h
Spike event reception register	0x05h
Target neuron population register	0x06h
Broadcast between all neighbour clusters	0x07h
Multicast between selected neighbour clusters op1	0x08h
Multicast between selected neighbour clusters op2	0x09h - 0x0Ch
Internal communication	0x0F
Debug mode	0x0Dh - 0x0Eh



On-chip Communication Protocols & Free Look-up Table Approach

ON-CHIP COMMUNICATION REGISTERS

Component	Memory Register Requirements [bit]					
	address	comm	Rx	Tx	timer	total
Cluster router	16	NA	NA	NA	NA	16
Tile router	20	NA	NA	NA	NA	20
Node router	4	4	10	30	14	62

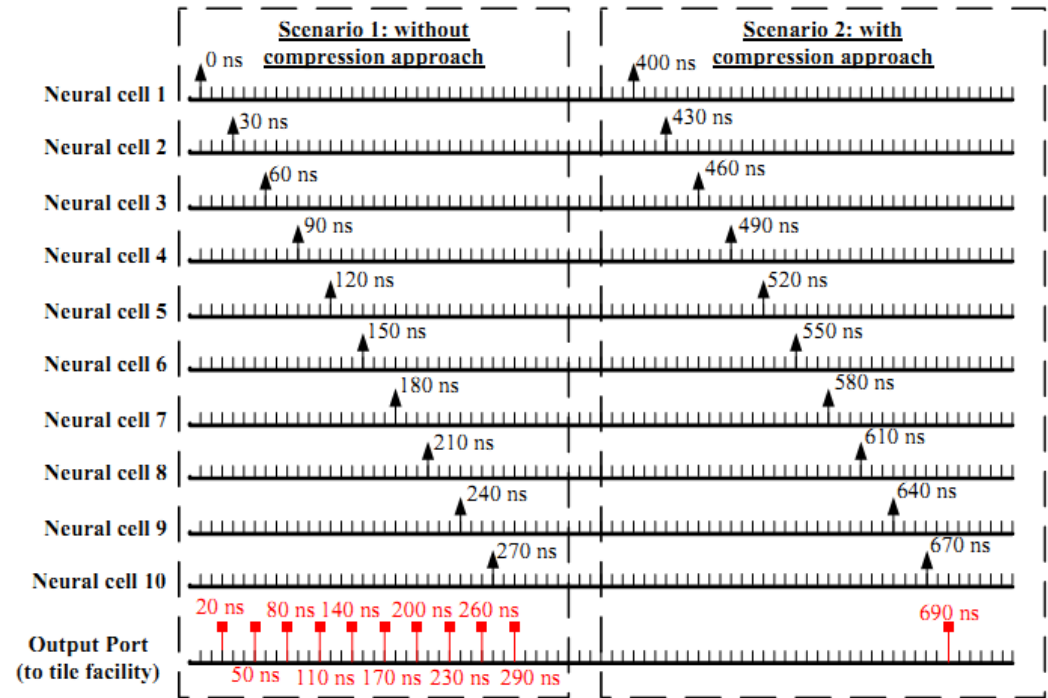
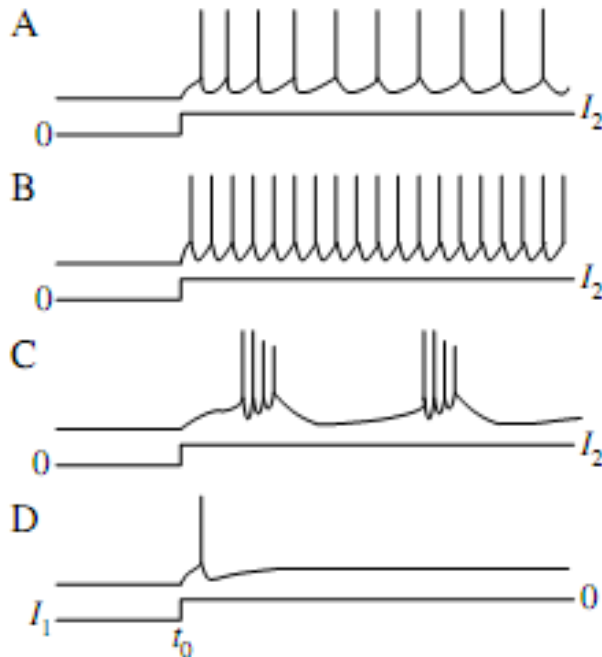
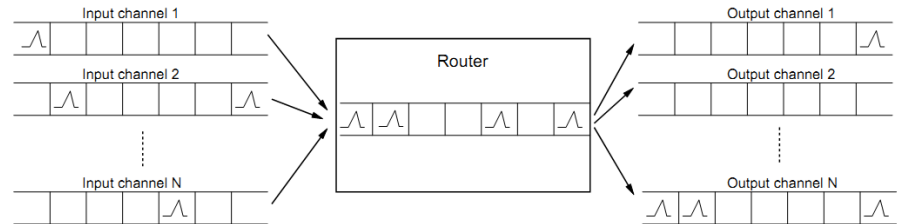
$$\begin{aligned} &+ \begin{aligned} &(1 \text{ cluster router}) \times (16 \text{ bits}) \\ &(4 \text{ tile routers}) \times (20 \text{ bits}) \\ &(40 \text{ node routers}) \times (62 \text{ bits}) \end{aligned} \\ &= \mathbf{2.576Kbit (400 \text{ neurons})} \end{aligned}$$

- The implemented approach shows a very significant reduction in memory size.
- Previous work shows memory requirements in the order of Mbits !!

Spike Event Compression Technique

Motivation:

- SNN traffic is slow ($ISI > 1\text{ms}$)
- Irregular pattern
- Polychronous Phenomena [Izhikevich'09]
- (i.e. More than 1 spike arriving at the same time)



Outline

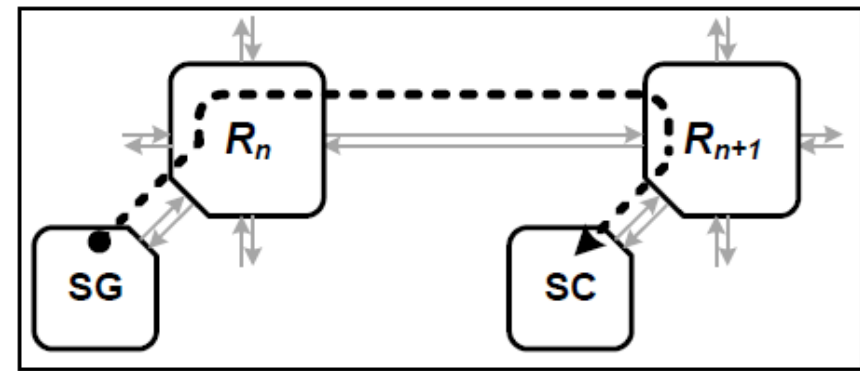
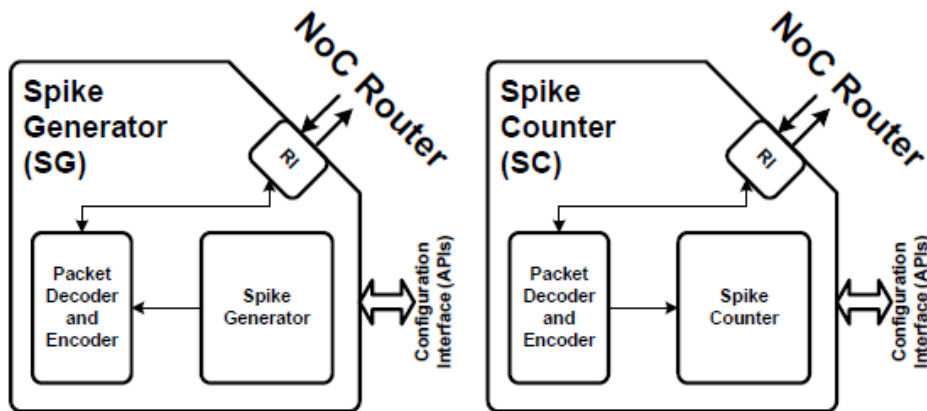
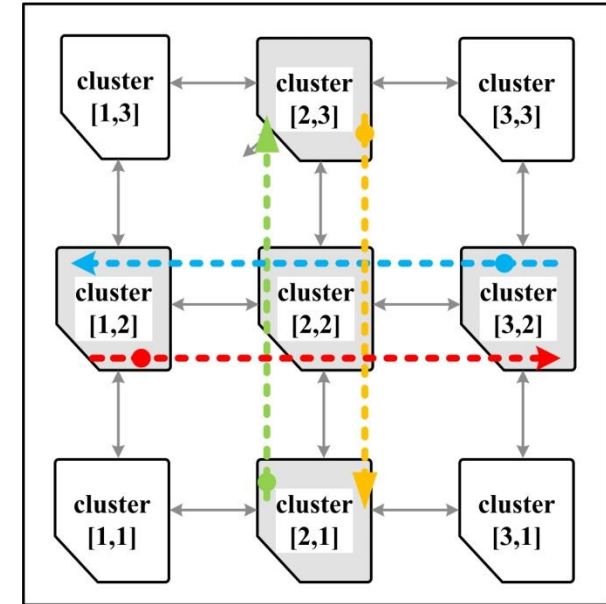
- Motivation and Challenges
- Hierarchical NoC EMBRACE Architecture
- Performance Analysis
- Take-home Message & Future Work



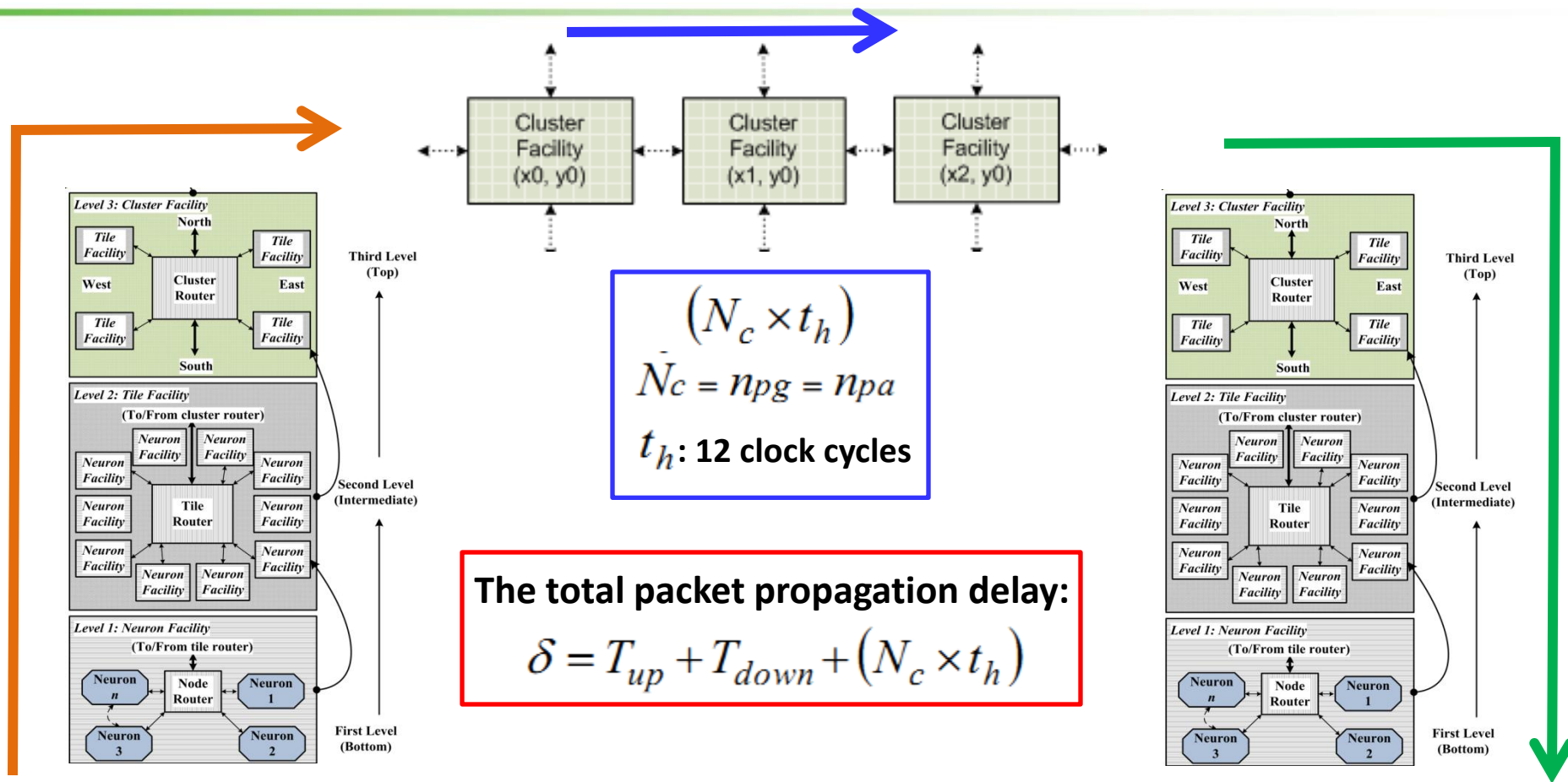
Experimental Setup

Methodology:

- **VHDL Simulation** of up to **50 x 50** array of clusters
- **FPGA implementation** of a 3x3 proof of concept array of clusters
- **100MHz** clock frequency per cluster & a **48-bits** packet
- **65-nm CMOS** technology (estimated)



Traffic Load Analysis



$$T_{up} = \tau_u + \left\lfloor t_u \times (n_{pg} - 1) \right\rfloor$$

$\tau_u : 24 \text{ clock cycles}$
 $t_u : 12 \text{ clock cycles}$

$$T_{down} = \tau_d + \left\lfloor t_d \times (n_{pa} - 1) \right\rfloor$$

$\tau_d : 30 \text{ cc}$
 $t_d : 1 \text{ cc}$

Traffic Load Analysis for Large Scale Scenarios

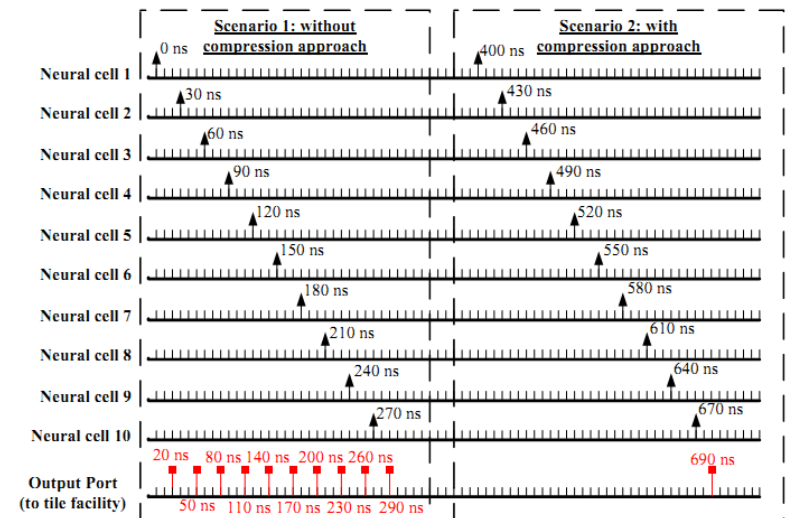
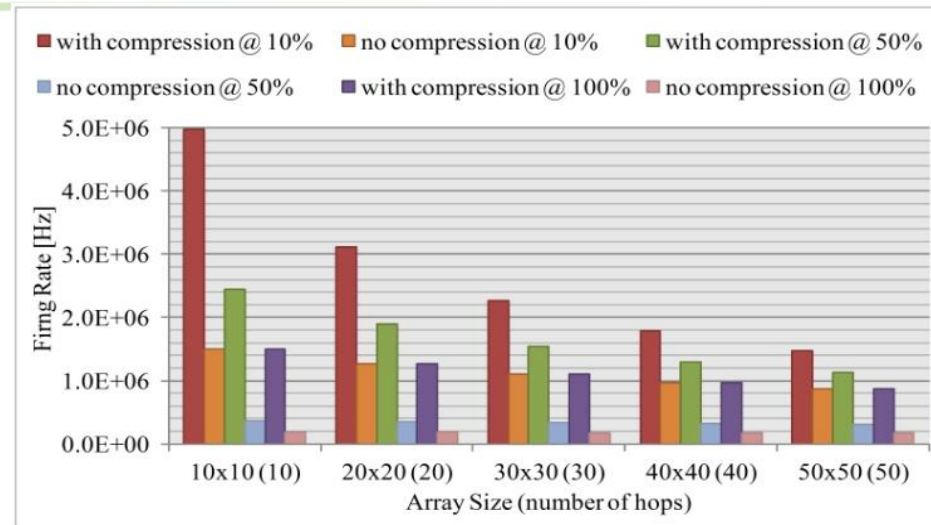
Typical biological spiking neurons show a firing rate around 100 Hz, but some others can show a firing rate up to 1KHz.

A maximum firing rate of ~5 MHz for a 10 hop scenario is highlighted using the compression approach.

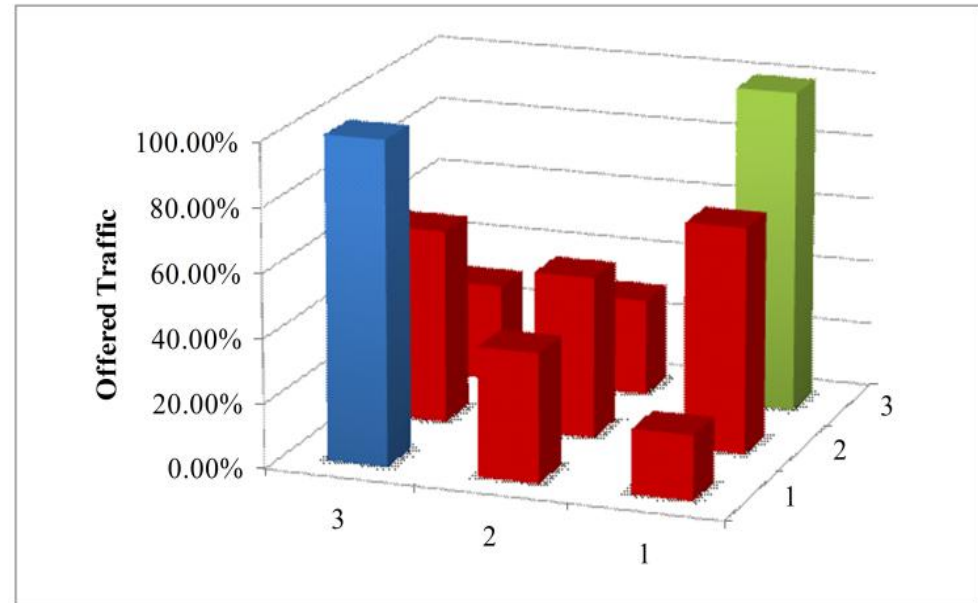
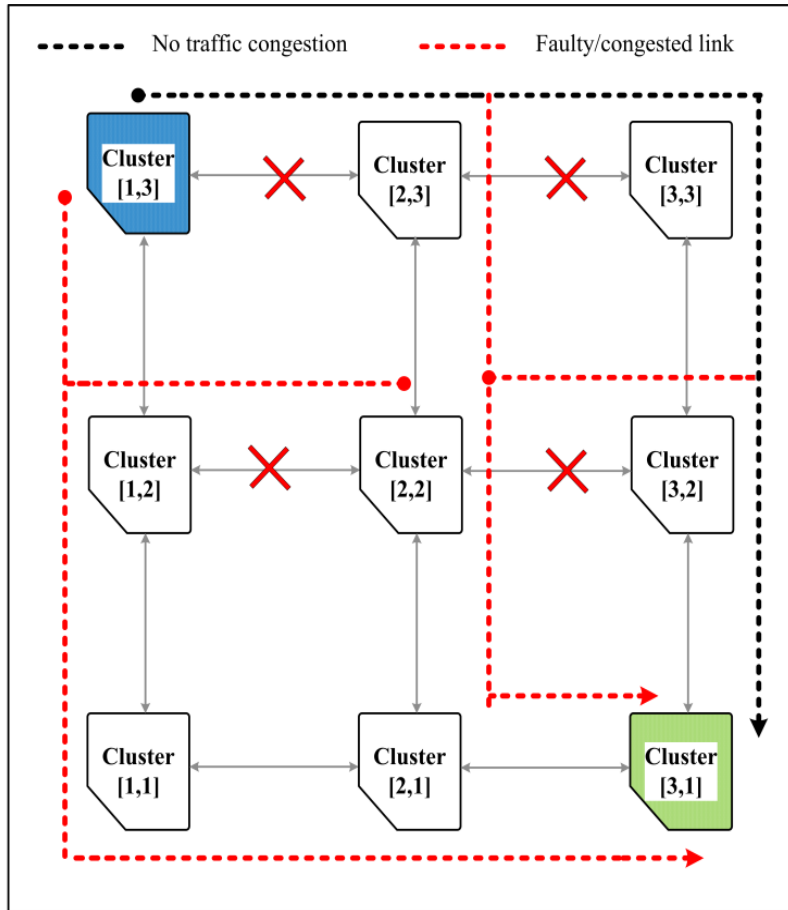
This offers a ~3.3x improvement compared to the same scenario without the compression technique.

In the 50 hop scenario, although the firing rate can decrease to 172 KHz when the compression technique is not used,

From a hardware point of view, if higher firing frequencies can be achieved, the platform can be used as a neural network hardware accelerator.



Adaptive Router Validation on FPGA

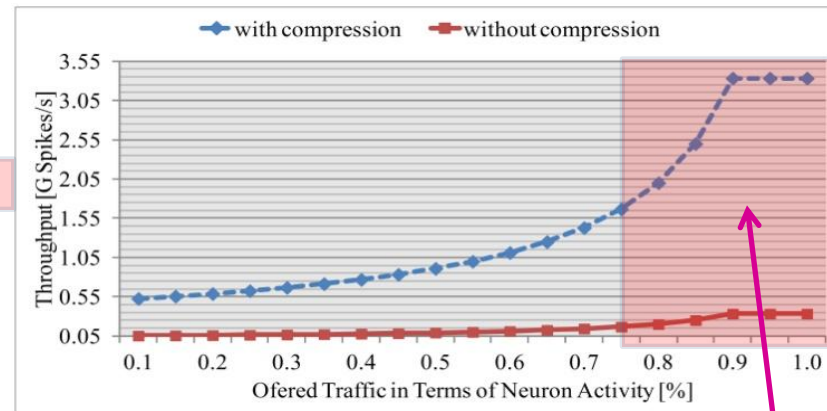


XY routing algorithm is used as a default routing mechanism when there is no traffic congestion

Throughput and Synthesis Results

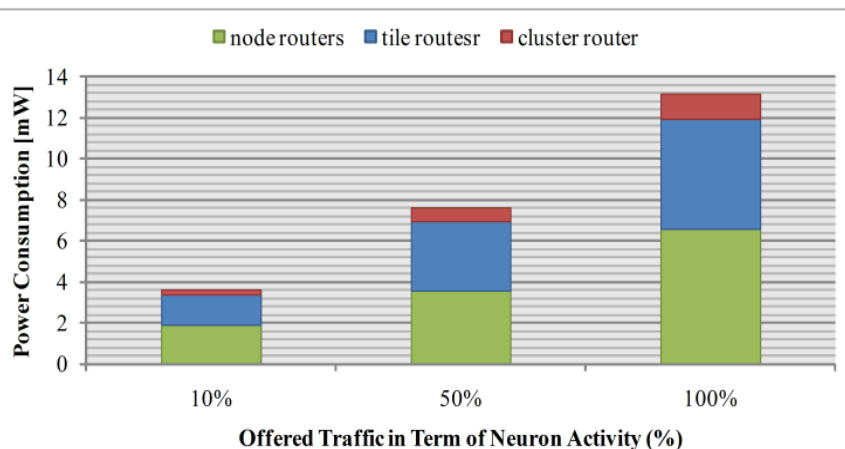
Area/power performance of router (65nm)

Component	Area [mm ²]			Power [mW]		
	Comb	seq	total	dynamic	static	total
Cluster	0.139	0.448	0.587	10.62	2.54	13.16
Cluster router	0.005	0.017	0.022	1.10	0.09	1.19
Tile router	0.055	0.015	0.070	3.86	0.27	4.13
Node router	0.002	0.005	0.007	0.41	0.03	0.44



Increased throughput under load testing

Power Consumption vs. Offered Traffic (65nm)



Project Reference	Quality of Service (QoS)	Congestion Mechanism	Throughput [Events/s]	Improvement
EMBRACE	Best Effort	Yes	3.33x10 ⁹	--
FACETS	Best/Guaranteed Effort	No	1.50x10 ⁹	2.2x
SpiNNaker	Best Effort	No	0.200x10 ⁹	16.5x

Proposed router outperforms existing approaches

Outline

- Motivation and Challenges
- Hierarchical NoC EMBRACE Architecture
- Performance Analysis
- Take-home Message & Future Work



Take-home Message & Future Work

- There are many problems associated with the development of efficient large scale SNN platform in hardware.
- A H-NoC approach is proposed as a way to overcome the intercommunication constrains currently experienced in the efficient realisation of SNNs in hardware.
- Future Work: Real-life SNN applications & Self-repair Mechanism based on the information received from the adaptive routing algorithm.



Acknowledgments

Snaider Carrillo Lindado is supported by
a Vice-Chancellor's Research Scholarship (**VCRS**)
from the University of Ulster

Thanks for your attention
Any question/feedback welcome

Snaider Carrillo

Email: carrillo_lindado-s@email.ulster.ac.uk

<http://isrc.ulster.ac.uk/Staff/SCarrillo/Contact.html>