



# Going for Brain-Scale Integration using FPGAs, TSVs and NOC- based Artificial Neural Networks- A Case Study

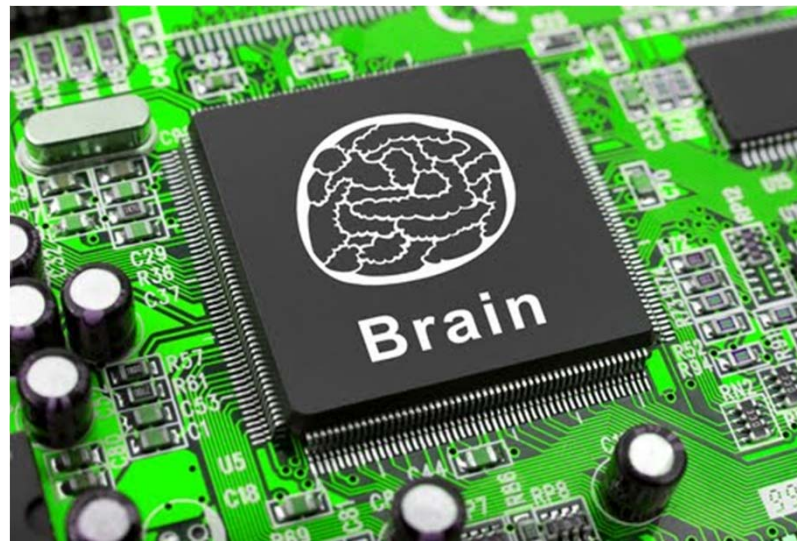
Nowshad Painda Mand, Johnny Öberg

Electronic System Dept  
School of information and Communication  
Technology  
Royal Institute of Technology(KTH), Sweden

# Outline

- Motivation
  - Related Work
  - Introduction
  - Neural Networks structure
  - Vision and Architecture
  - Experimental Results
  - Conclusion and Future Work
-

# Let's build a brain...



# Motivation-Tasks

(Biological systems do with great speed)

- Common sense reasoning
  - Simple tasks like walking/talking
  - Imprecise info interpretation
  - Generalization/Learning
  - Voice, image, Video processing etc
-

# Motivation

## Biological Systems

Biological Systems	
Processor	<ul style="list-style-type: none"><li>• Simple</li><li>• Low speed</li><li>• Large Number</li></ul>
Memory	<ul style="list-style-type: none"><li>• Distributed</li><li>• Content Addressable</li></ul>
Computation	<ul style="list-style-type: none"><li>• Distributed</li><li>• Parallel</li><li>• Self-learning</li></ul>
Power Consum	<ul style="list-style-type: none"><li>• Low</li></ul>
Reliability	<ul style="list-style-type: none"><li>• Very Robust</li></ul>
Complex Arith	<ul style="list-style-type: none"><li>• No</li></ul>

# Motivation-Realizations

- Better understanding of brain
    - Individual cell
    - Part of brain
    - Malfunction/Diseases/Better drugs
  - Artificial Intelligence
    - Pattern recognition
    - Speech and image processing
    - Robotics
  - New computer architecture
    - Based on massive parallelism
      - For problem not solvable by current architectures
-

# Related Work-i

- SpiNNaker (Manchester Uni)
    - Tiny ARM968 CPUs on Chip, On chip/ inter-chip network, 65,536-18(1 Million Cores), (Partially Completed, 2013).
  - C2S2 (IBM + 5 Uni, 2008)
    - First SW simulator now turned HW
    - Based on crossbar and not scalable
  - FACET and BraiscaleS (19 European Uni)
    - Analog/digital Hybrid, Limited neurons can be simulated
-

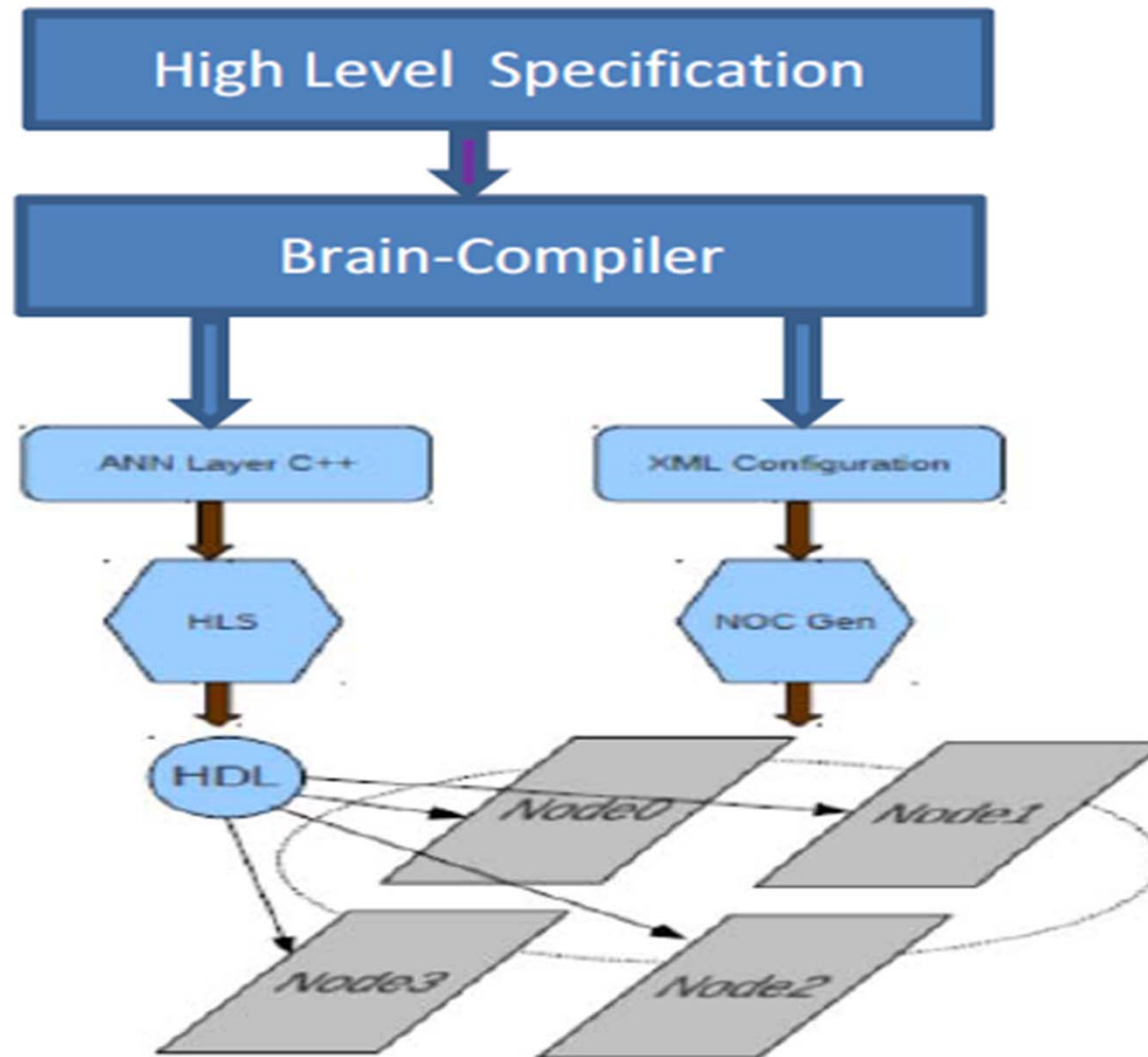
# Related Work-ii

- BioRC (Southern California)
    - Synthtic Cortex using nano tubes and Graphene
    - Started very recently
  - IBM Brain-Inspired Computer
    - One million Neurons
    - 256 million synapses
    - 5.4BT ,2nd largest CMOS Chip in the World
    - 4096 distributed cores On-Chip Mesh network
    - Over 400 million bits of local on chip memory, around 100kbits of memory per core
-

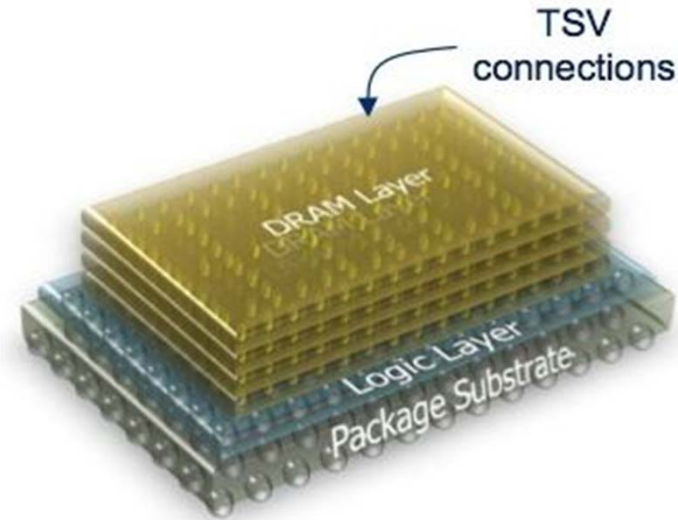
# Automated Tool for Experimental Setup

- FPGA base system
  - Plasticity through reconfigurability
  - Best neural structure
  - Experimental platform
-

# Our Vision



# Architecture



- BPU's on chip
- BPU Board containing BPU Chips
- Packet-switch NOC for BPU's
- Off-chip protocol for BPU Boards
- For sufficient bandwidth chip-stack of DRAMs on BPU chips

---

\*Figure from Micron Technologies- IBM 3D Chips

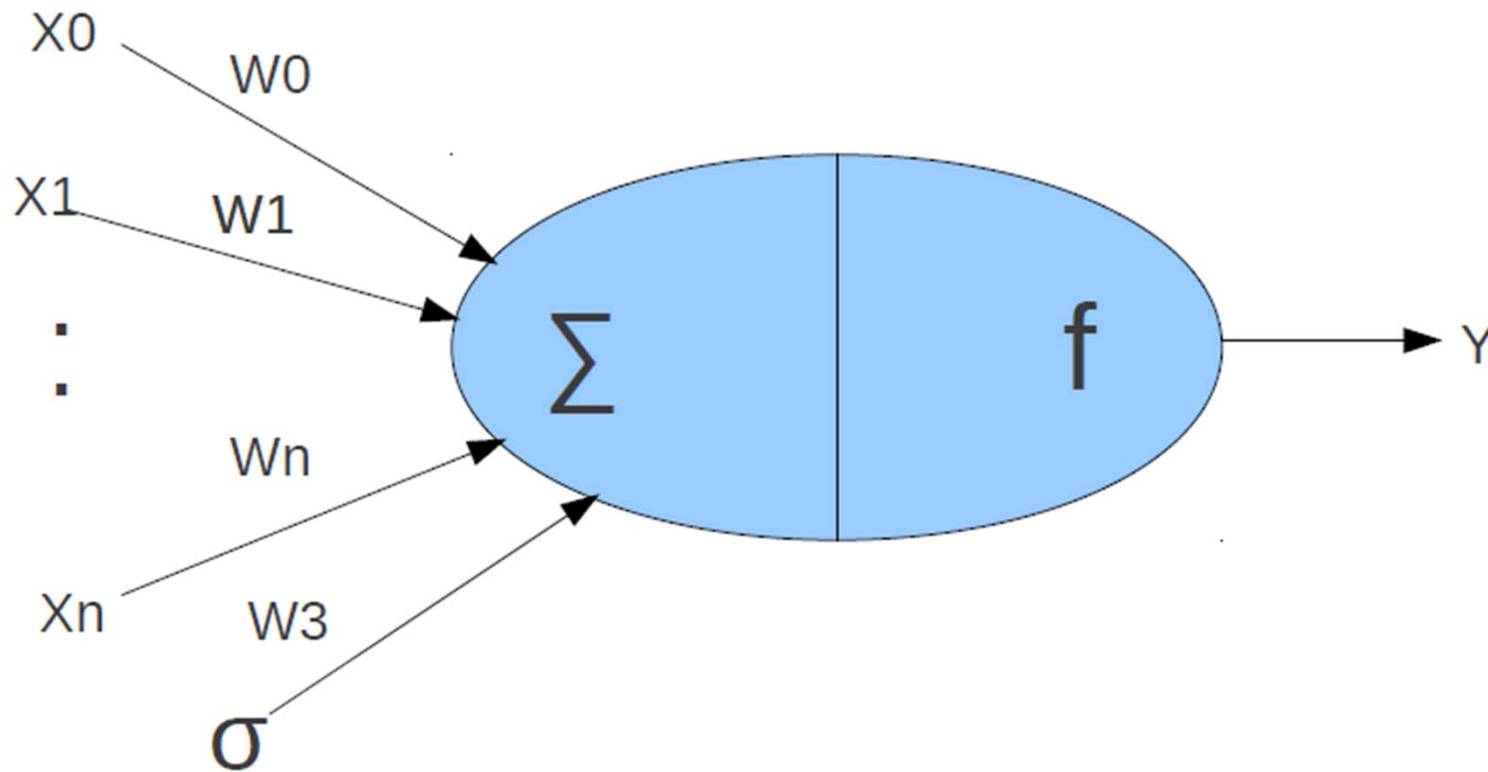
# Introduction-i

- ANN modelled as
    - Software, Hardware or Hybrid
  - Software realisation
    - Sequential and not suitable for real time apps
  - Hardware realization
    - Ideally suited for real time
    - Conventional HW does not scale up
    - Non linear inter-neuron communication
    - Plasticity, fan-in/fan-out & power consumption
-

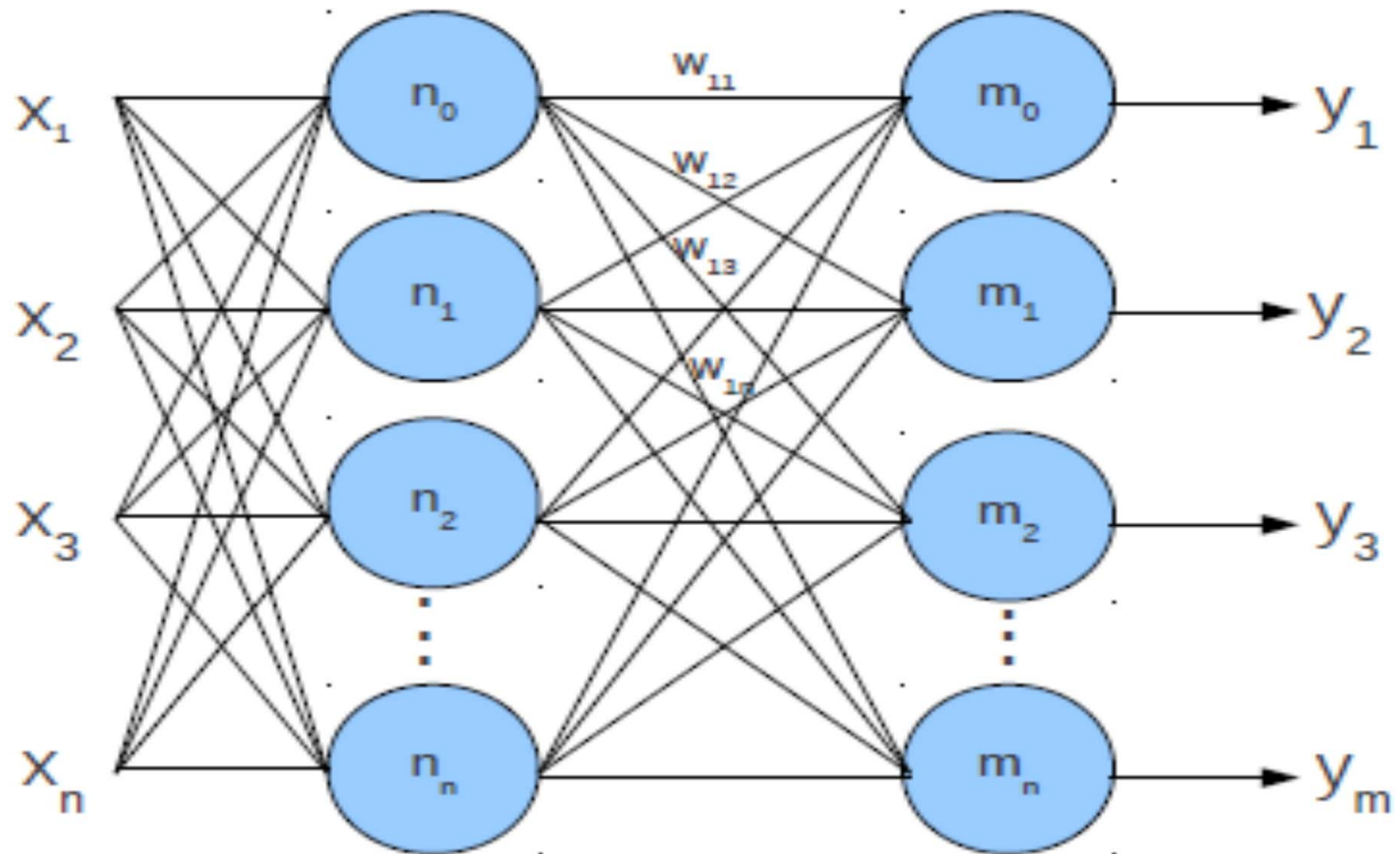
# How many FPGAs to Build a Brain ?

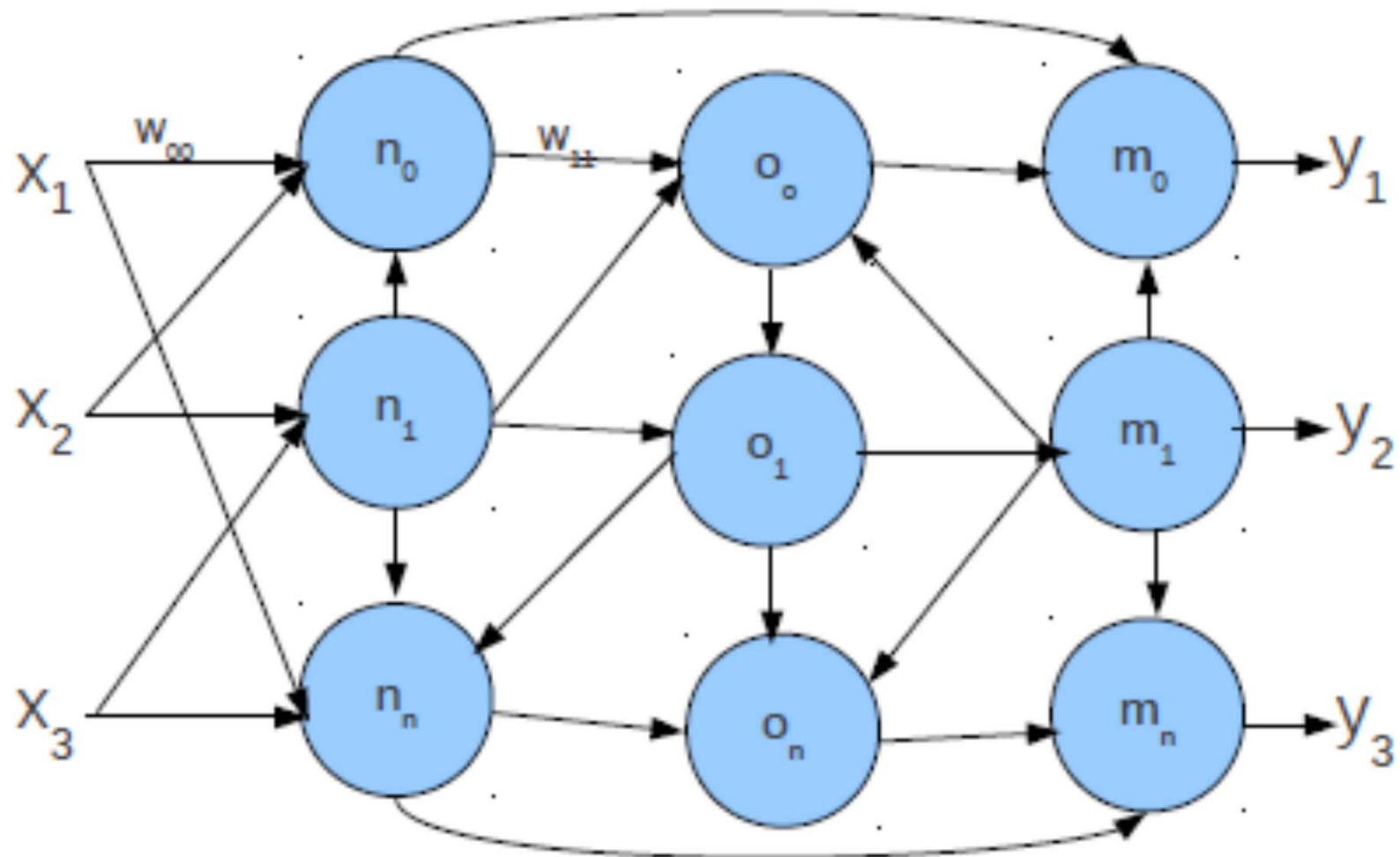
- High level neural net model
  - Design-space exploration
  - Synthesized model using HLS tools
  - Mapping of synthesize model on NOC based FPGA platform
-

# Basic Neuron



# Neural Network-i





# Requirements

- On Avg 1000 weights & multiply with 1000 for every neuron firing
  - $10^{11}$  neurons & 1% firing every 10ms
  - Need  $10^3 \times 10^{11}$  weighted memory
  - Need to fetch  $10^3 \times 10^{11} \times 1\% / 10_{\text{ms}} = 10^{14}$  input weights per sec
-

# Design-Space Exploration

- No of BPUs depends on
    - Speed of the DRAM-don't shrink with technology
    - Amount of DRAM on chip-stack-more TSVs means simple logic (row/column addr), pipelining to speedup increase area so speed won't be an issue
    - No of TSVs possible
    - Storage limited by no of dies in stack(Technology) & no of TSVs
    - Finally no of TSVs limiting factor
-

# Design-Space Exploration

Global Level, W2W, D2W or D2D 3D- stacking	2009-2012	2012-2015
Minimum TSV diameter	4-8 $\mu\text{m}$	2-4 $\mu\text{m}$
Minimum TSV pitch	8-16 $\mu\text{m}$	4-8 $\mu\text{m}$

-Stacking TSVs diameter as per ITRS Roadmap,2012

---

# Design-Space Exploration

- Conservative approach with TSV pitch  $8\mu\text{m}$ 
    - TSV density  $1\text{cm}^2/64\mu\text{m}^2 = 1,562,500 \text{ TSV}/\text{cm}^2$
    - Never need  $\log_2(10^{15}) = 50$  bit to address whole memory
    - One calculation need 1000 input weight
    - Address bits and result bits are tiny and can be ignored
-

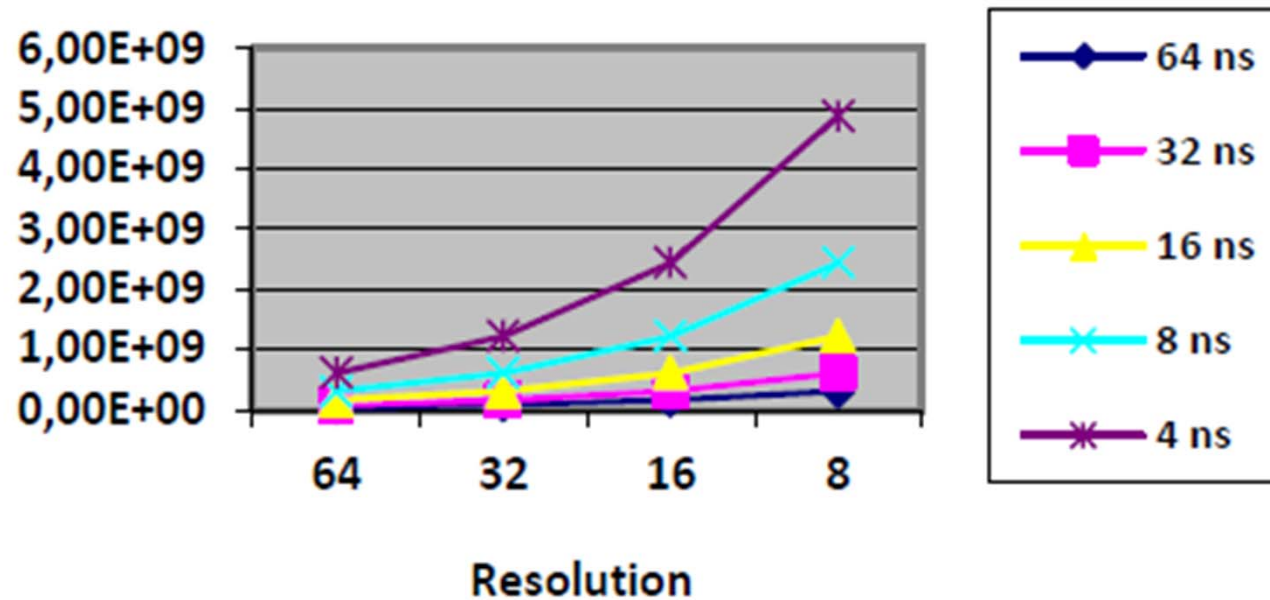
# Design-Space Exploration

$$\begin{aligned} & \text{NrNeuronCalcsPerCycleCm2} \\ & \quad \text{TSV}_{\text{density}} \\ = & \frac{\quad}{\text{resolution} * \text{NrInputWeightPairsPerNeuron}} \end{aligned}$$

---

# Design-Space Exploration

# Neuron Calculations per cm<sup>2</sup>



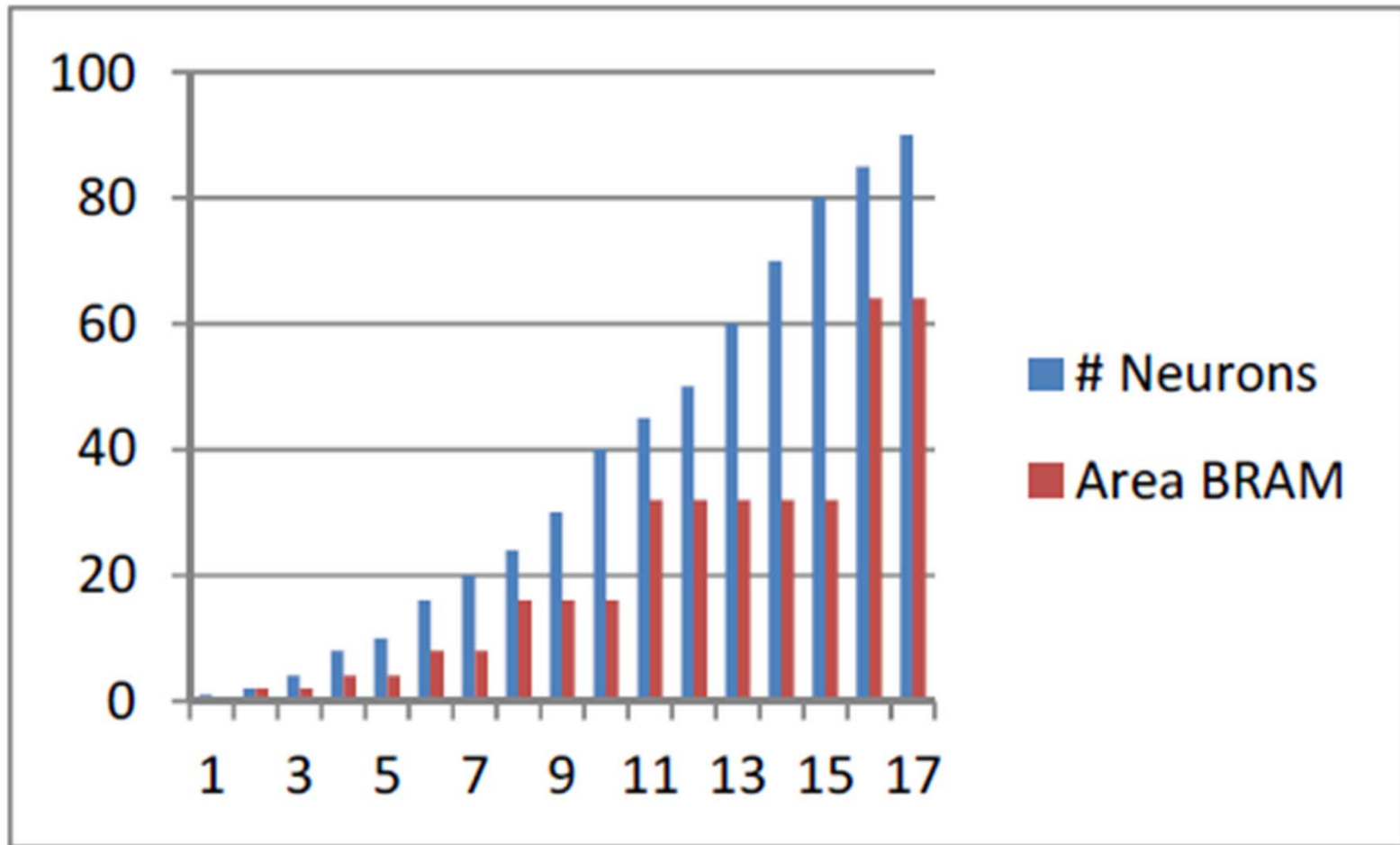
# Design-Space Exploration

Input-Weight Resolution [bits]		64	32	16	8	4
DRAM Size [Byte/brain]		$8 \cdot 10^{15}$	$4 \cdot 10^{15}$	$2 \cdot 10^{15}$	$10^{15}$	$5 \cdot 10^{14}$
#Chips	100	$8 \cdot 10^{13}$	$4 \cdot 10^{13}$	$2 \cdot 10^{13}$	$10^{13}$	$5 \cdot 10^{12}$
	1,000	$8 \cdot 10^{12}$	$4 \cdot 10^{12}$	$2 \cdot 10^{12}$	$10^{12}$	$5 \cdot 10^{11}$
	10,000	$8 \cdot 10^{11}$	$4 \cdot 10^{11}$	$2 \cdot 10^{11}$	$10^{11}$	$5 \cdot 10^{10}$
	100,000	$8 \cdot 10^{10}$	$4 \cdot 10^{10}$	$2 \cdot 10^{10}$	$10^{10}$	$5 \cdot 10^9$

# Experimental Results

- For Vivado HLS two cases are considered
    - Case-1 weights are modelled as being read from input ports
    - Case-2 weights are modelled as being synthesized as on-chip BRAM
  - NOC synthesized separately
    - NOC node consumed around 750LUTs
    - One BPU around 18000LUTs
-

# No. of Neurons vs FPGA Area



# Weight Data Type Dependencies

Weight Data Type	Area [LUT]	Delay [clk-cc]
64-bit Double Serial	13807	5563
64-bit Integer Serial	15862	5005
64-bit Integer fully parallel	64000 <sup>1</sup>	1
32-bit Single Serial	7201	5563
32-bit Integer Serial	8957	5004
16-bit Integer Serial	5555	5004
8-bit Integer Serial	3504	5004
4-bit Integer Serial	2652	5003

---

<sup>1</sup>Vivado could not handle fully parallel case, assumed manual calc

# Throughput Considerations

<b>Input-Weight Resolution [bits]</b>	<b>64</b>	<b>32</b>	<b>16</b>	<b>8</b>
<b>Production rate #Neuron Calculations/cycle/chip</b>	<b>54.9</b>	<b>109</b>	<b>220</b>	<b>439,4</b>
<b>Consumption rate #Neuron Calculations/cycle/chip</b>	<b>15</b>	<b>31</b>	<b>62</b>	<b>125</b>
<b>Fraction of required rate</b>	<b>27%</b>	<b>28%</b>	<b>28%</b>	<b>28%</b>

# Most Feasible Design

DRAM Access Time [ns]	Resolution (#bits per Input-weight pair)	#Chips ( $2.25 \text{ cm}^2$ )	DRAM Size (Bytes per chip)
8	64	$5.5 \cdot 10^5$	$1.5 \cdot 10^9$
	32	$2.7 \cdot 10^5$	
	16	$1.3 \cdot 10^5$	
	8	$6.6 \cdot 10^4$	

# Conclusion

- Design will be area limited for FPGAs
  - Using DRAMs ANN models can be mapped to modern FPGAs
  - Equivalent of full brain will need ~66,000 chips each with 1.5Gb DRAM stack
-

# Future work

- More types of neuron models
  - Other types of activation functions
  - Build specialized NOC for ANN
  - Build complete model with limited neurons
  - HLS tool for massive parallelism
-

# Questions??

---

# Design-Space exploration

DRAM Access Time [ns]	Resolution (#bits per Input-weight pair)	#Neuron Calculations per second per cm <sup>2</sup>	Area needed for entire brain [cm <sup>2</sup> ]	#Chips (2.25 cm <sup>2</sup> )
64	64	$3.8 \cdot 10^8$	$26 \cdot 10^5$	$12 \cdot 10^5$
	8	$30.5 \cdot 10^8$	$3.3 \cdot 10^5$	$1.5 \cdot 10^5$
32	64	$7.6 \cdot 10^8$	$13 \cdot 10^5$	$5.8 \cdot 10^5$
	8	$61 \cdot 10^8$	$1.6 \cdot 10^5$	$7.3 \cdot 10^4$
16	64	$15.3 \cdot 10^8$	$6.6 \cdot 10^5$	$2.9 \cdot 10^5$
	8	$122 \cdot 10^8$	$8.2 \cdot 10^5$	$3.6 \cdot 10^4$
8	64	$30.5 \cdot 10^8$	$3.3 \cdot 10^5$	$1.5 \cdot 10^5$
	8	$244 \cdot 10^8$	$4.1 \cdot 10^4$	$1.8 \cdot 10^4$
4	64	$61 \cdot 10^8$	$1.6 \cdot 10^5$	$7.3 \cdot 10^4$
	8	$488 \cdot 10^8$	$2.0 \cdot 10^4$	$9.1 \cdot 10^3$
2	64	$122 \cdot 10^8$	$8.2 \cdot 10^4$	$3.6 \cdot 10^4$
	8	$977 \cdot 10^8$	$1.0 \cdot 10^4$	$4.6 \cdot 10^3$
1	64	$244 \cdot 10^8$	$4.1 \cdot 10^4$	$1.8 \cdot 10^4$
	8	$1,950 \cdot 10^8$	5,120	2,276

# Design-space exploration

- Cycle time is determined by DRAM access time
  - More aggressive pitch 4um
    - We get values 4 times larger and 4 times smaller in terms of area
  - Assume DRAM access time of 1ns for the next table
-

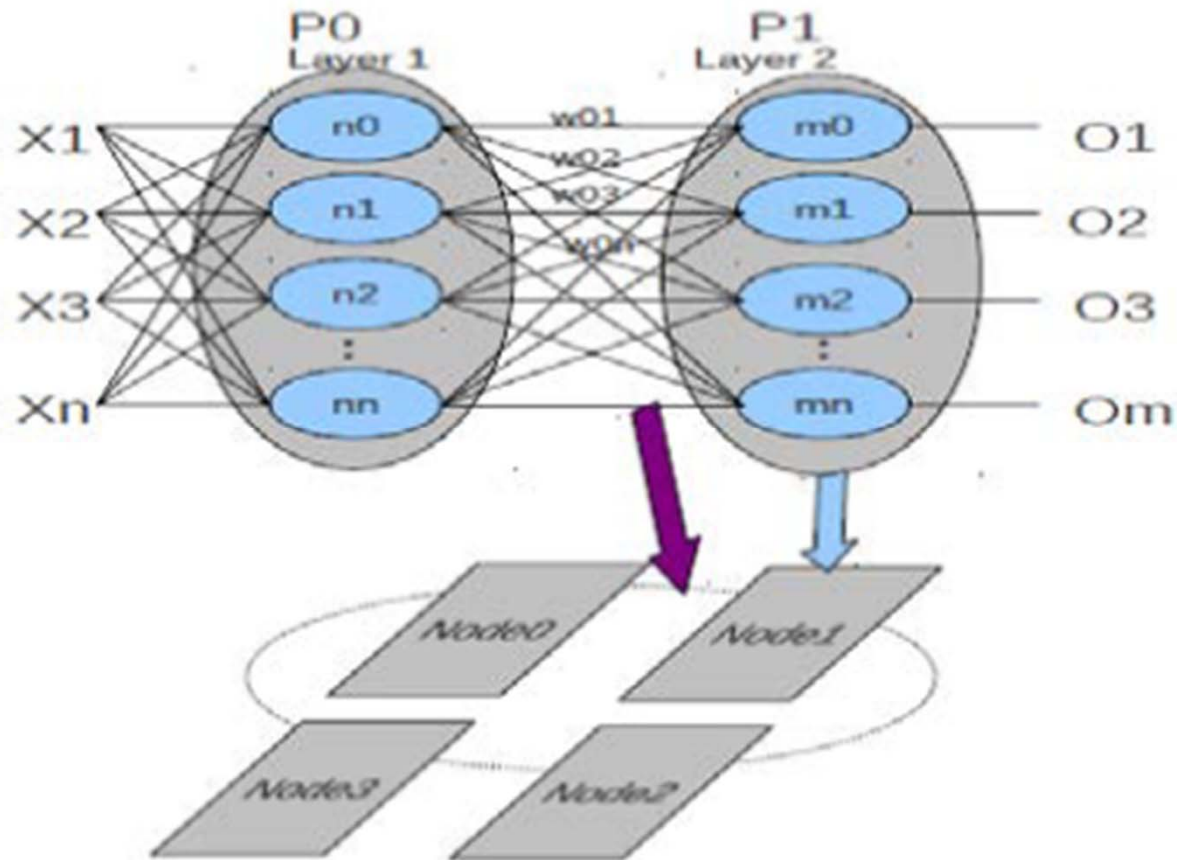
# Design-Space exploration

<b>Input-Weight Resolution [bits]</b>	64	32	16	8
<b># Neuron Calculations/cycle/cm<sup>2</sup></b>	97.7	195.3	390.6	781.3
<b># Chips</b>	4,551	2,276	1,138	569
<b># BPU Boards (100 chips/board)</b>	45.6	22.8	11.4	5.7

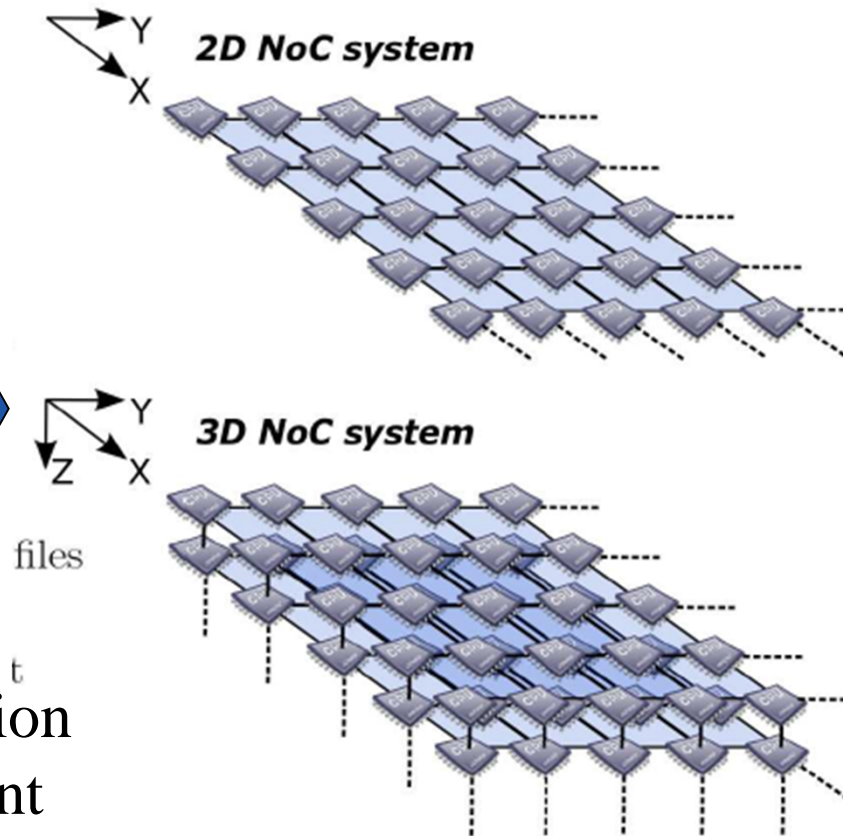
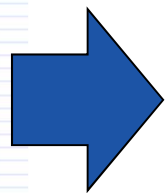
# BPU Implementation

- $1000 \times \text{prec} \times \text{res}[\text{bits}]$  need to be fetched for single Neuron calc
  - For actual BPUs we selected multi-layer perceptron model
  - Minicolumn & macrocolumns
  - Packet switch NOC to connect BPUs
-

# Mapping BPUs to NOC platform



# NOC Generator



Sopc/mhs files Creation  
Partial Pin Assignment



# Vivado HLS tool

- Xilinx state of the art HLS tool
  - ANN's layer behavioral model in c++
  - Single layer is synthesized
  - Leaving low-level optimizations to tool
  - For new point in design-space just change any parameter and re-run
-

# Nr of Neuron vs Resource requirement & latency

# Neu	BRAM (18K)	DSP48E	FF	LUT	Lat (clk-cc)
1	0	14	13591	15871	616
2	0	28	13610	16251	668
4	0	28	13619	16336	768
8	0	28	13628	16421	968
10	0	28	13628	16770	1068
20	0	28	13637	16892	1568
30	0	28	13642	17043	2068
40	0	28	13646	17013	2568
50	0	28	13652	17180	3068
60	0	28	13652	17170	3568
70	0	28	13655	17254	4068
80	0	28	13655	17135	4568
90	0	28	13655	17319	5068
100	0	28	13662	17311	5568