# FPGA Neural Networks

Hardware Neural Net Entertainment System (HNES)
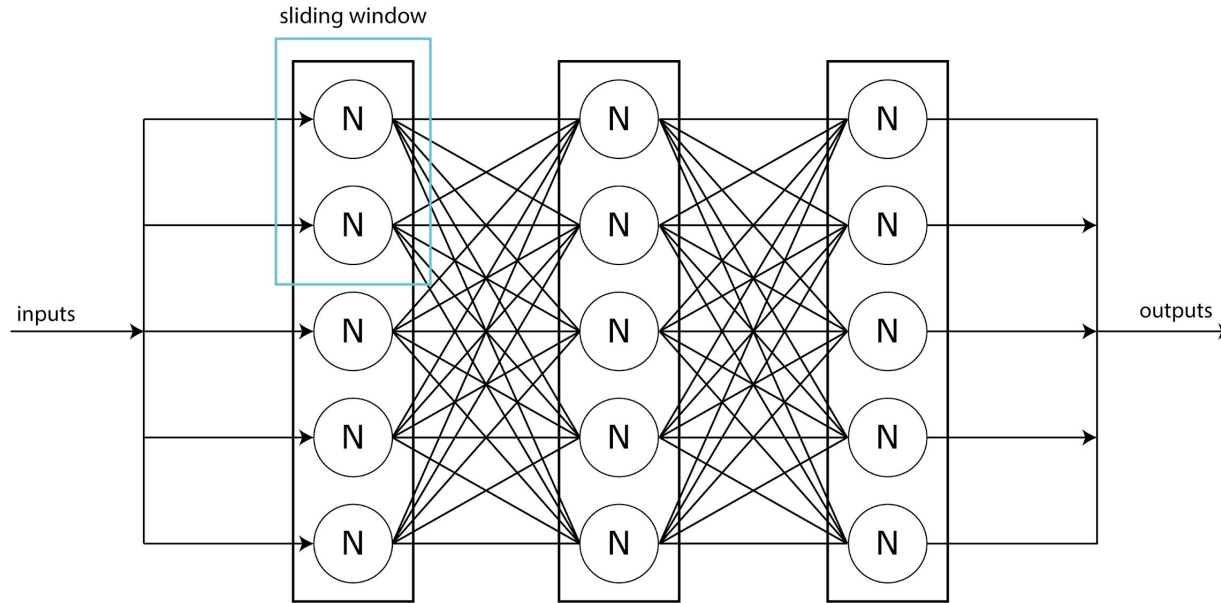
Sebastian Bartlett, Josh Noel

# Project Overview

Motivation: Applications of ML on FPGA is an active area of research

Inference on Feed-Forward network topology

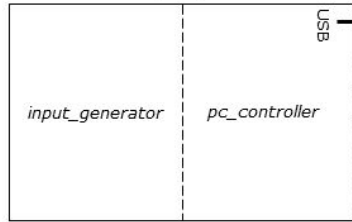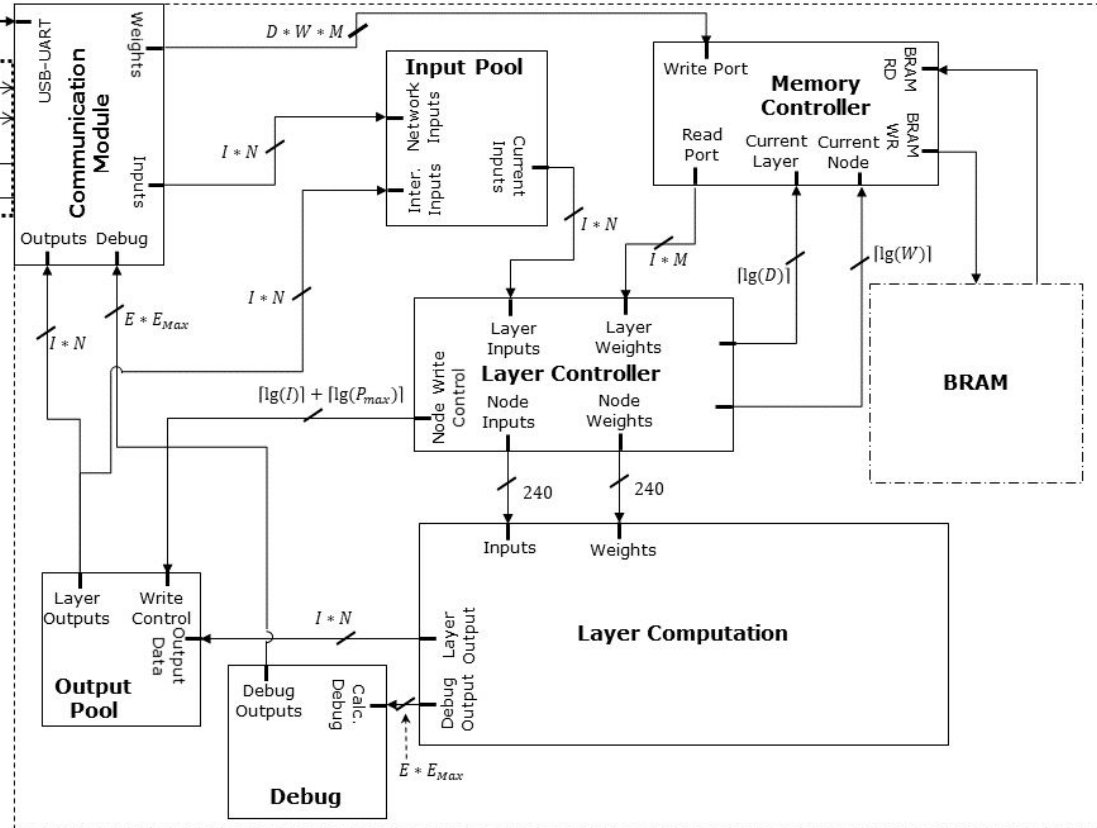Interface to PC for network inputs & outputs

sliding window

inputs

outputs

Network Architecture
(conceptual, not module)

Parallelism => P : Length ( Sliding window )

Network Depth => D : Length ( Layers )

Network  Width => W : Length ( Nodes Per Layer )

**PC**

**Nexys 4 DDR**

## Procedure

1. **Comm. Module** - Input Transfer
   a. PC transfers weights & network inputs

2. **Layer Controller** reads **Input Pool** and weights from **Mem Ctrl**. Ensures all layer outputs in **Output Pool**

3. **Output Pool** from layer *i - 1* is **Input Pool** for layer *i*

4. Repeat 2 & 3 until **Output Pool** contains network outputs

5. **Comm. Module -** Output Transfer
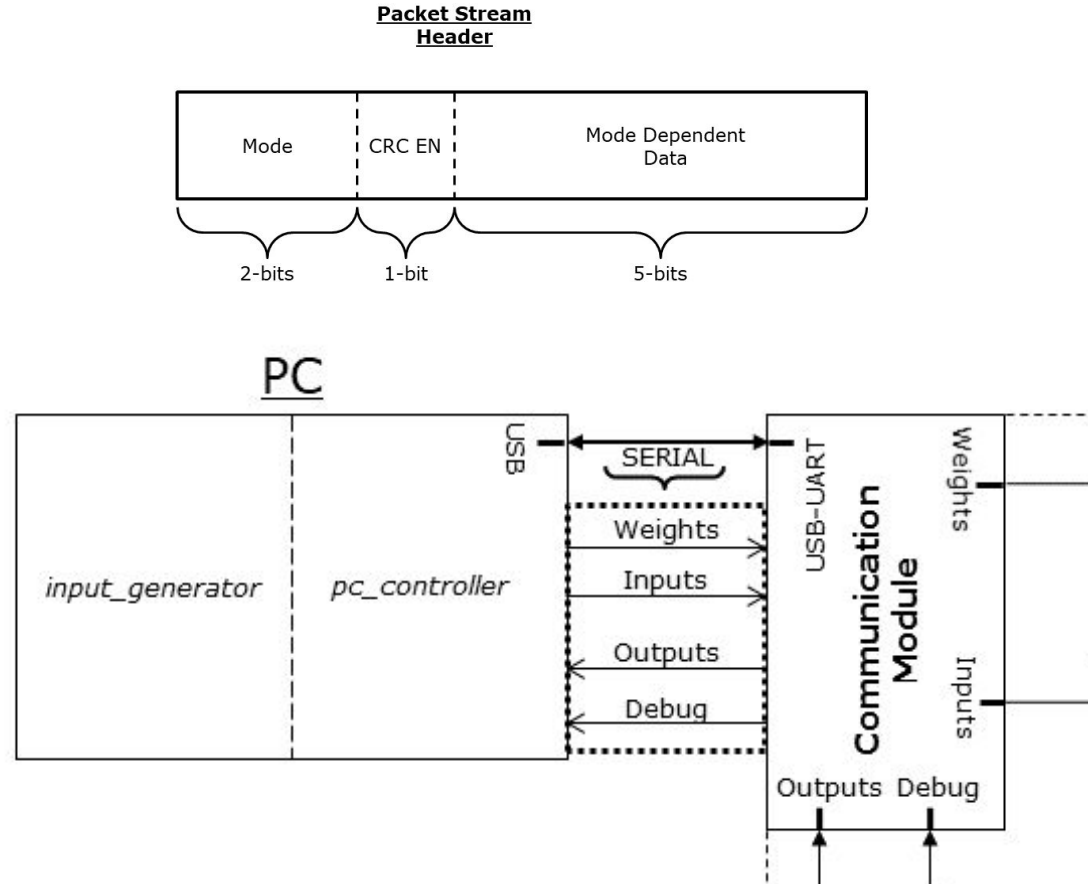
# Implementation Overview

# Communication Module

Communication over USB-UART - 10 bit data packets

**Packet Stream:** [Header, 9 data packets]

Header defines data packet type

Must divide data across multiple packets

| Mode | CRC EN | Mode Dependent Data |
|------|--------|---------------------|
| 2-bits | 1-bit | 5-bits |



PC

input_generator | pc_controller

USB — SERIAL — USB-UART

Weights
Inputs
Outputs
Debug

Communication Module
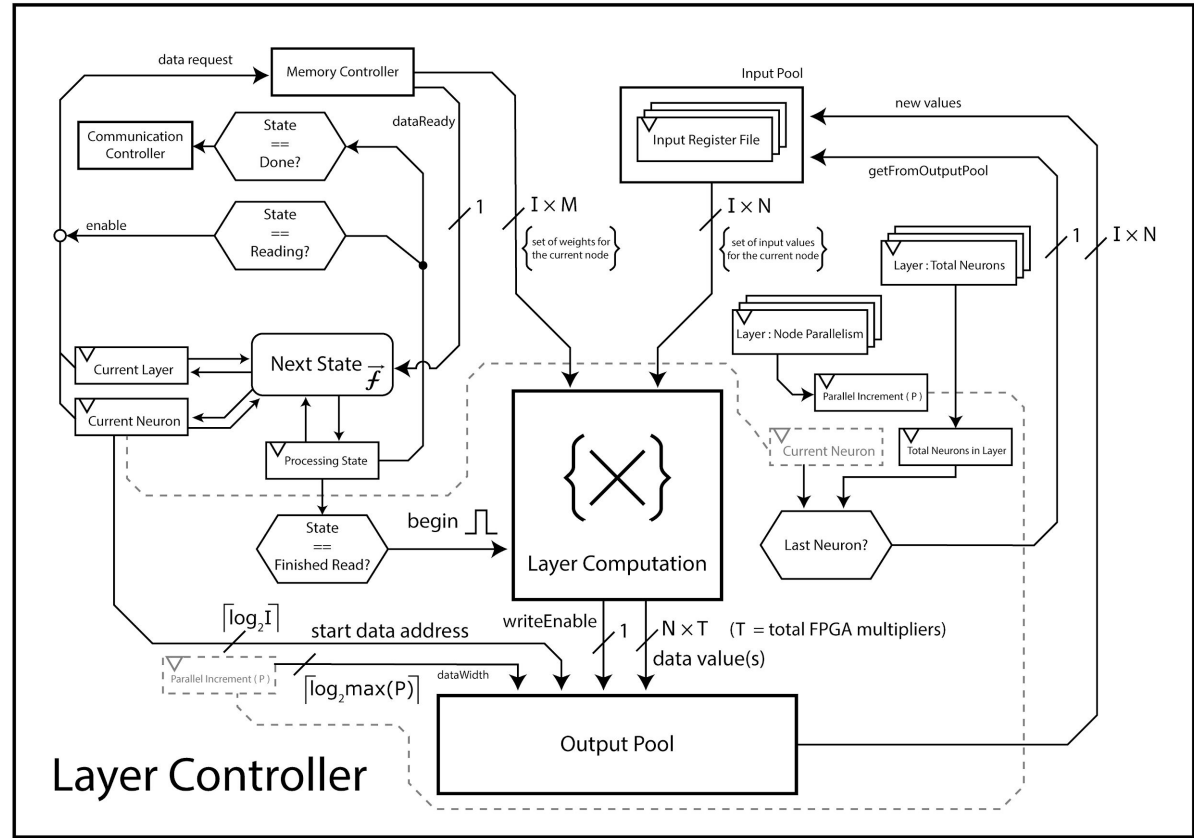
Weights

Inputs

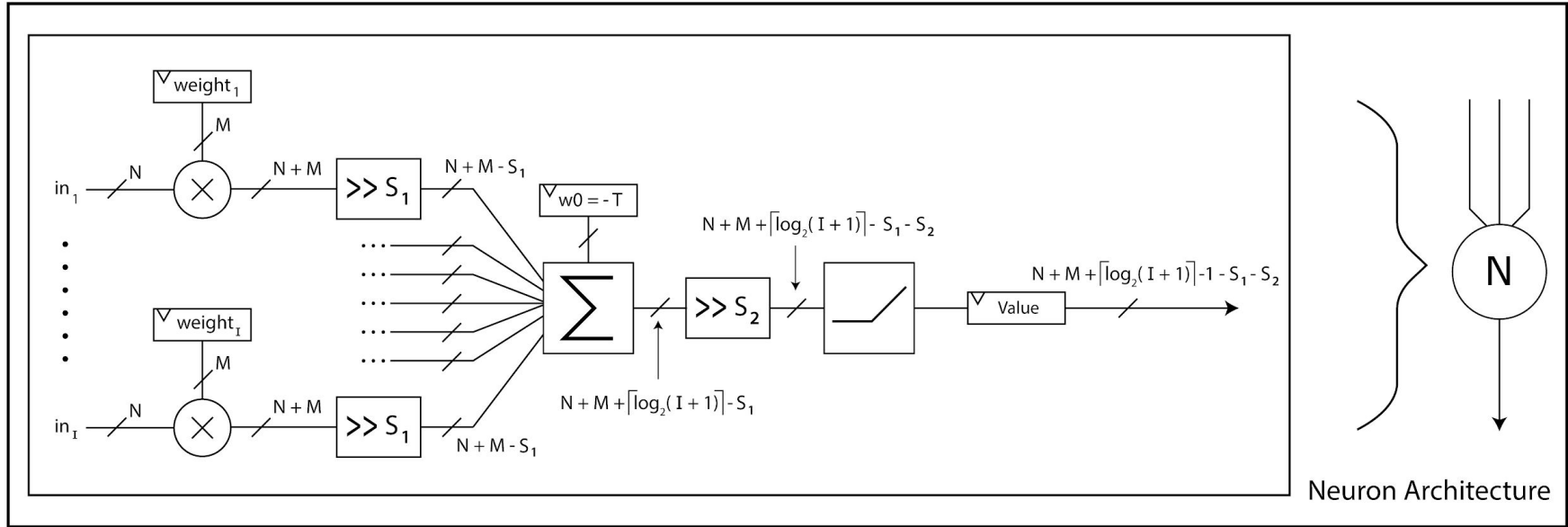Outputs   Debug

# Layer Controller

Controls all state transfers

Routes input pool & weight memory reads to **Layer Computation**

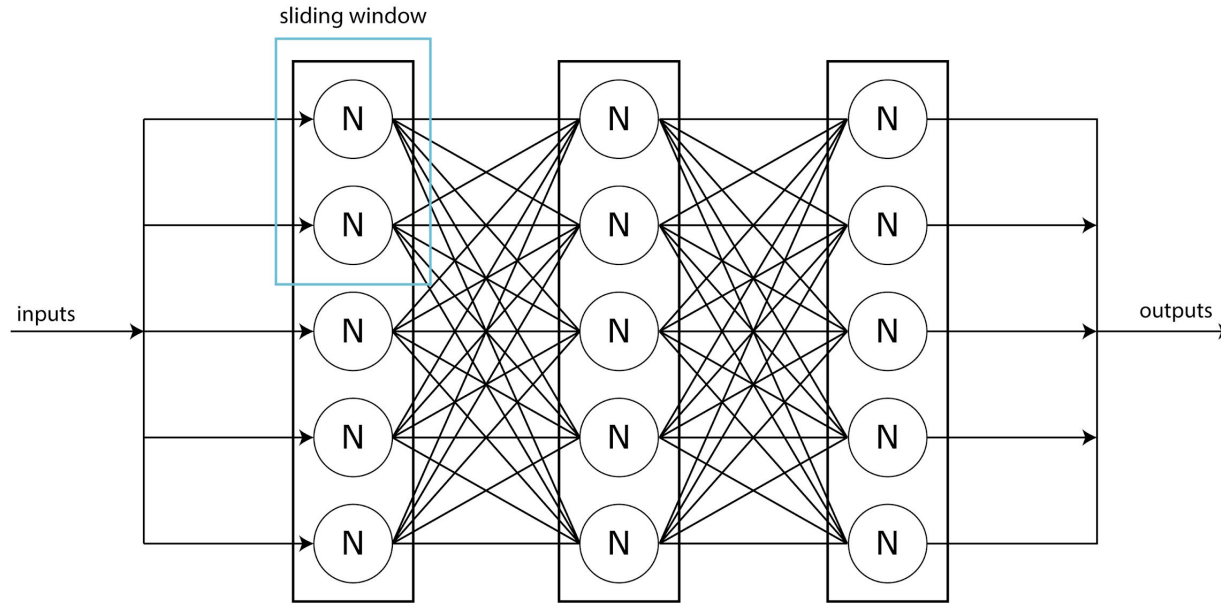Routes **Layer Computation** output to correct addresses in **Output Pool**

# Layer Computation Module



Neuron Architecture

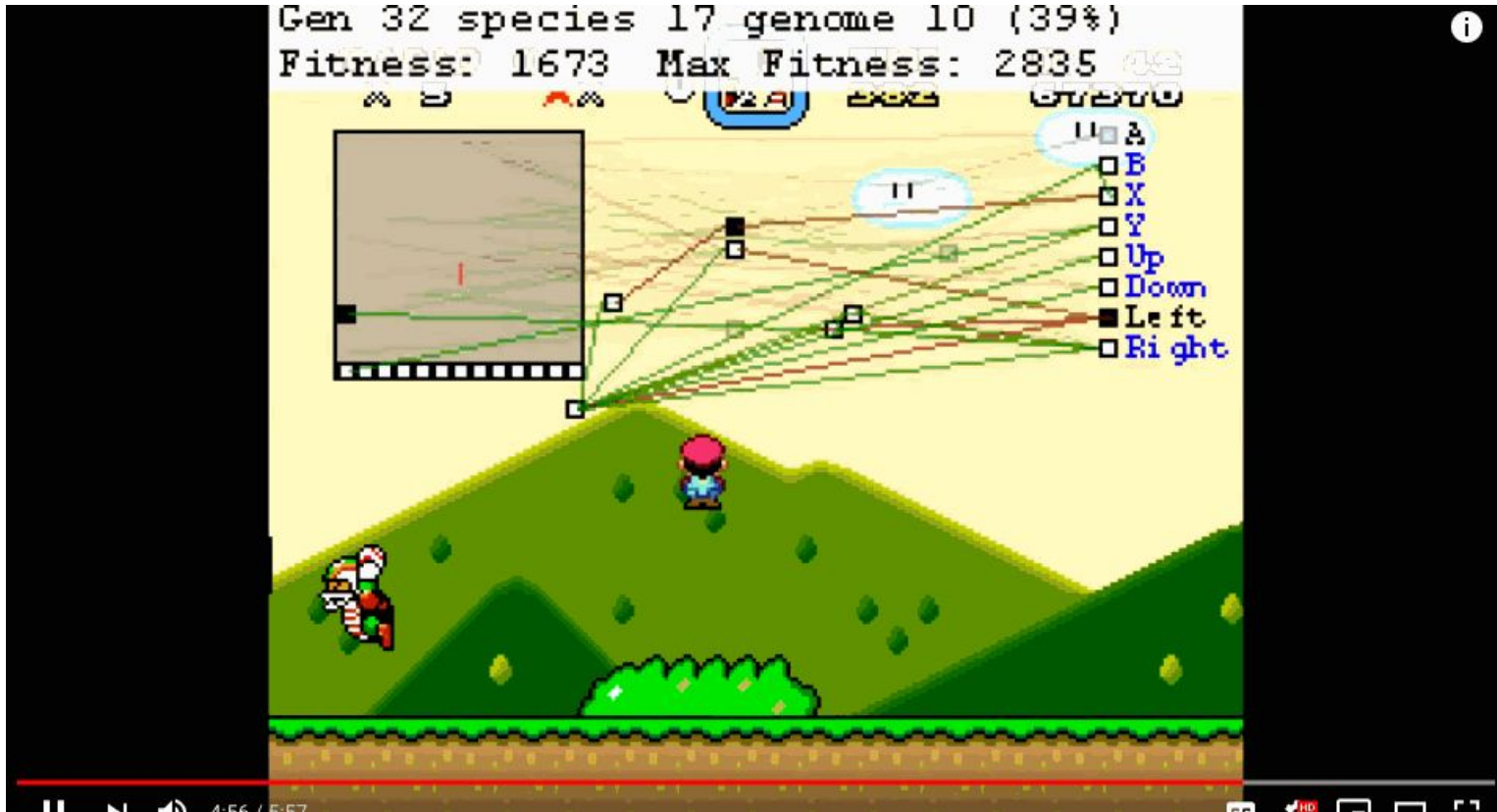# of neurons calculated at once (parallelism) is bounded by DSP slices

sliding window

inputs

outputs

Network Architecture
(conceptual, not module)

Parallelism => P : Length ( Sliding window )

Network Depth => D : Length ( Layers )

Network Width => W : Length ( Nodes Per Layer )
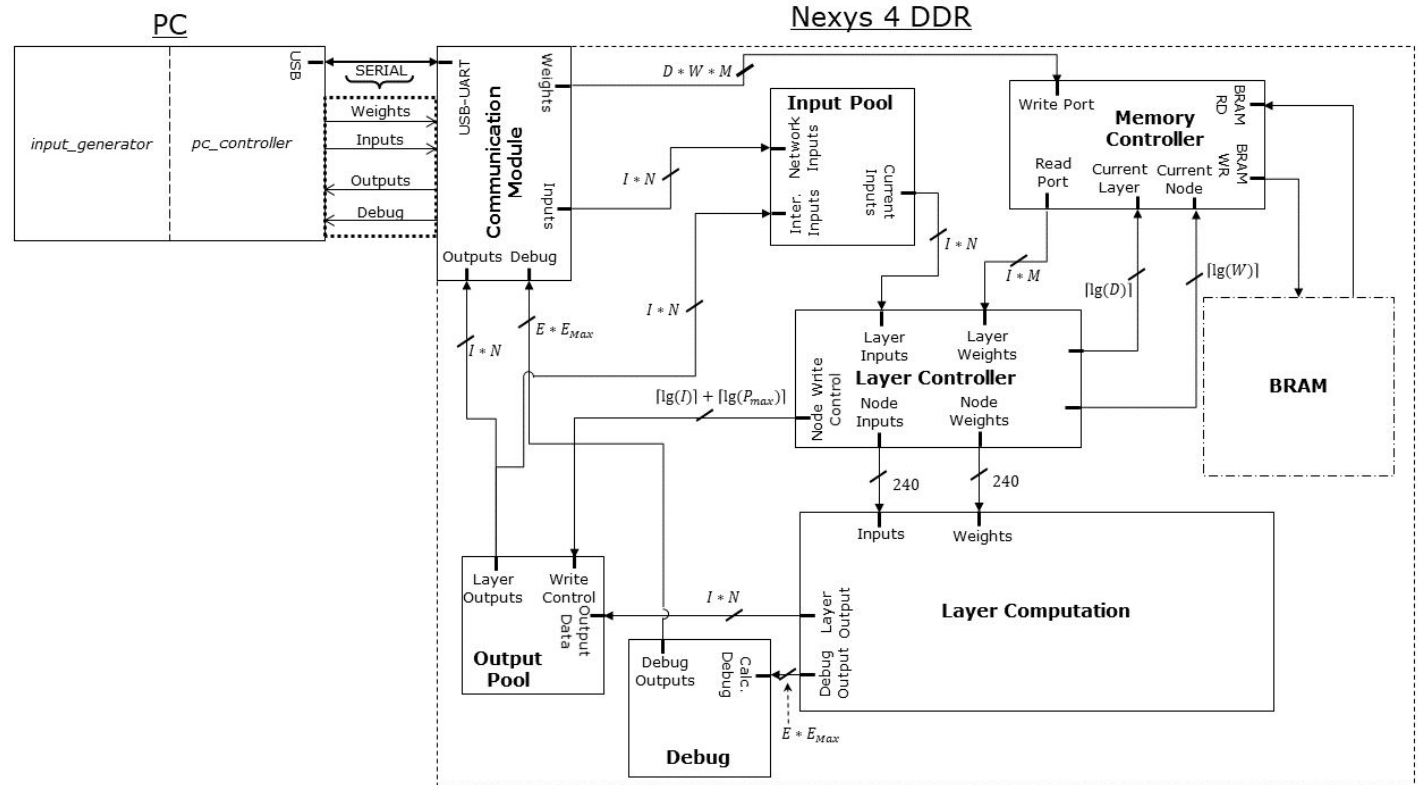
# Timeline

- 11/16 - Module Implementation
- 11/30 - Module Debugging, Integration Testing
- 12/10 - Neural Network Application and Implementation
  - Base Goals
    - Calculate single layer outputs based off PC inputs
      - FPGA behaves as a hardware matrix multiplier
    - Calculate output layer of multilayer feed-forward network
      - This requires calculating intermediate (hidden) layers
      - This enables inference to be performed
      - More complex topologies can be interfaced on PC-side
  - Stretch Goal
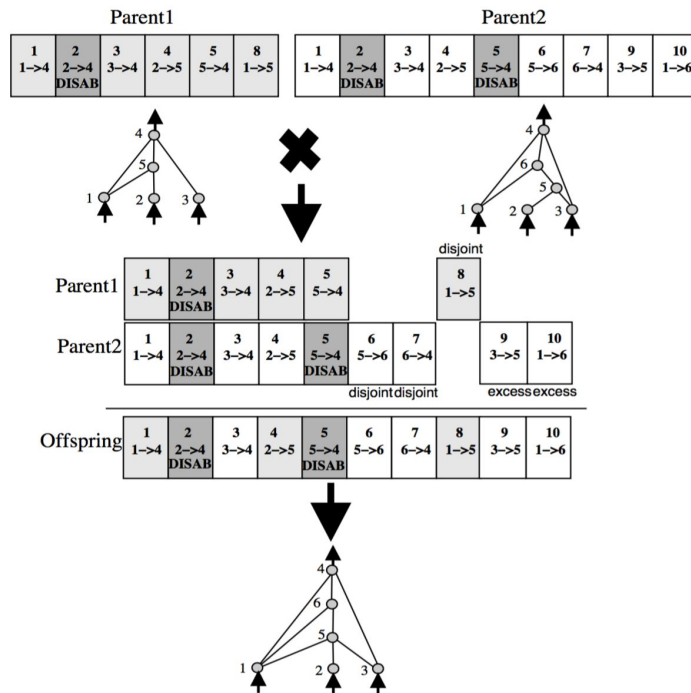    - Interface with NES emulator to play Super Mario

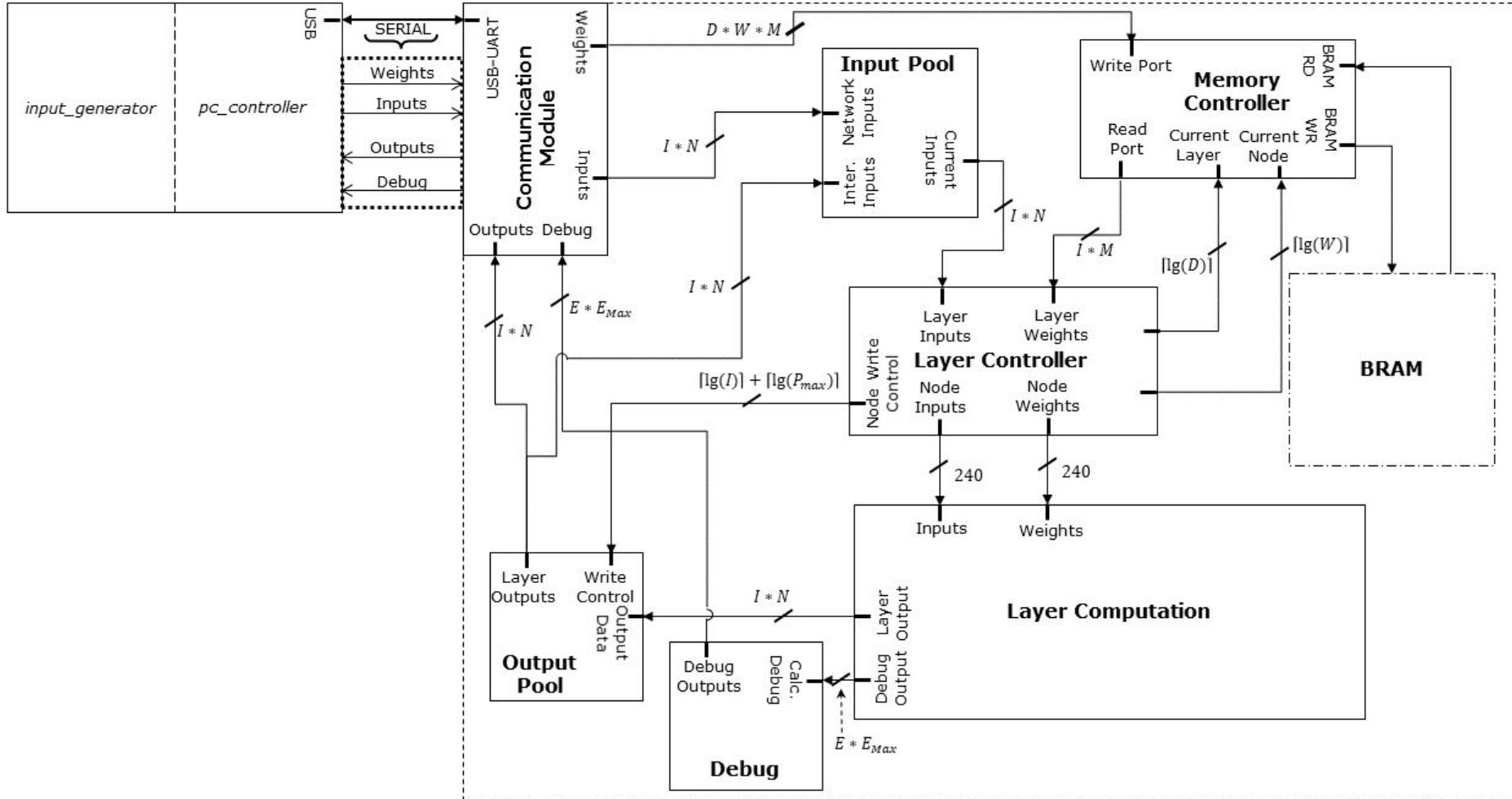Mar I/O by "SethBling" (YouTube)

# Questions?

(Reference Slides)

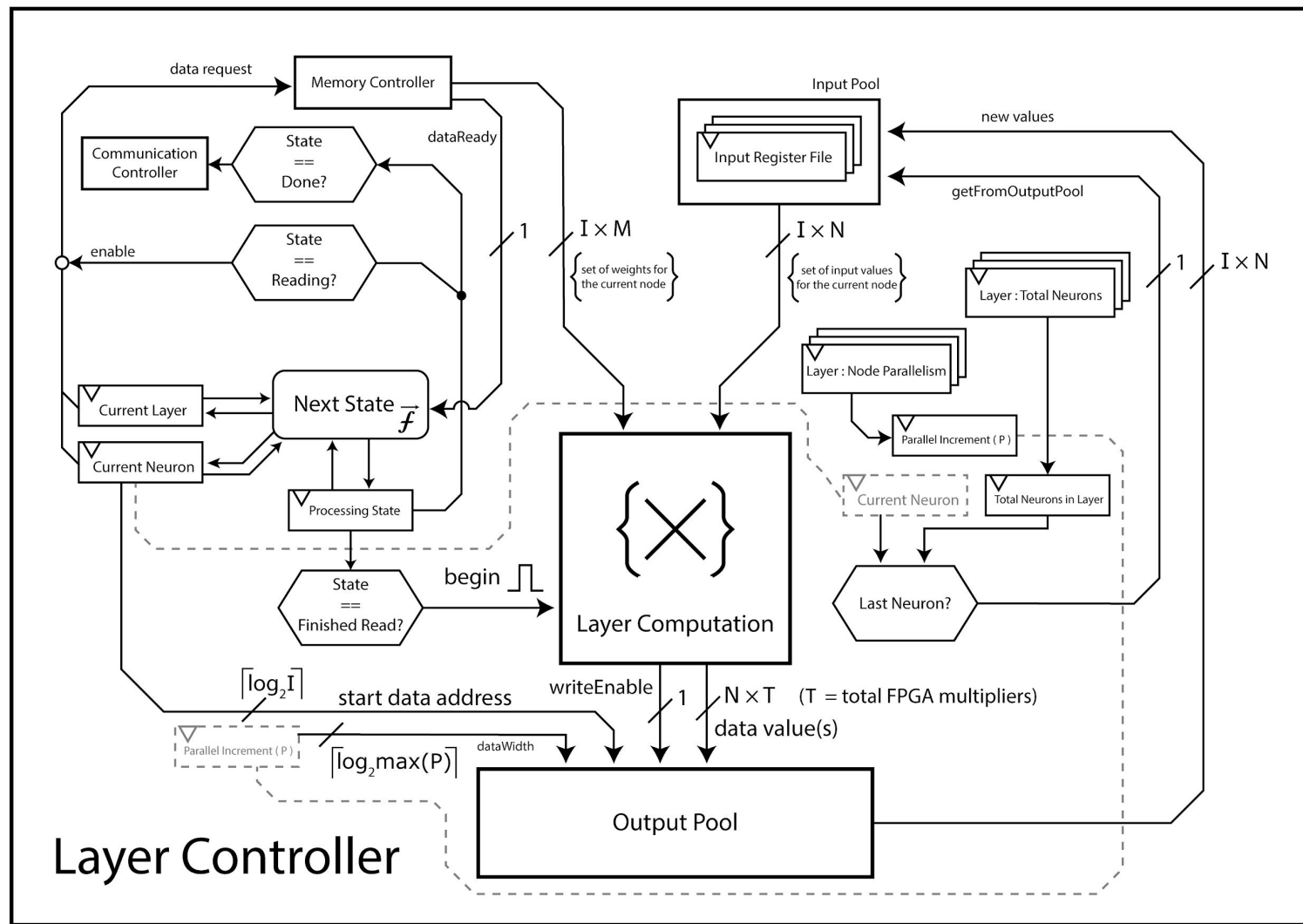# Source: Evolving Neural Networks through Augmenting Topologies

**Kenneth O. Stanley** kstanley@cs.utexas.edu Department of Computer Sciences, The University of Texas at Austin, Austin, TX 78712, USA
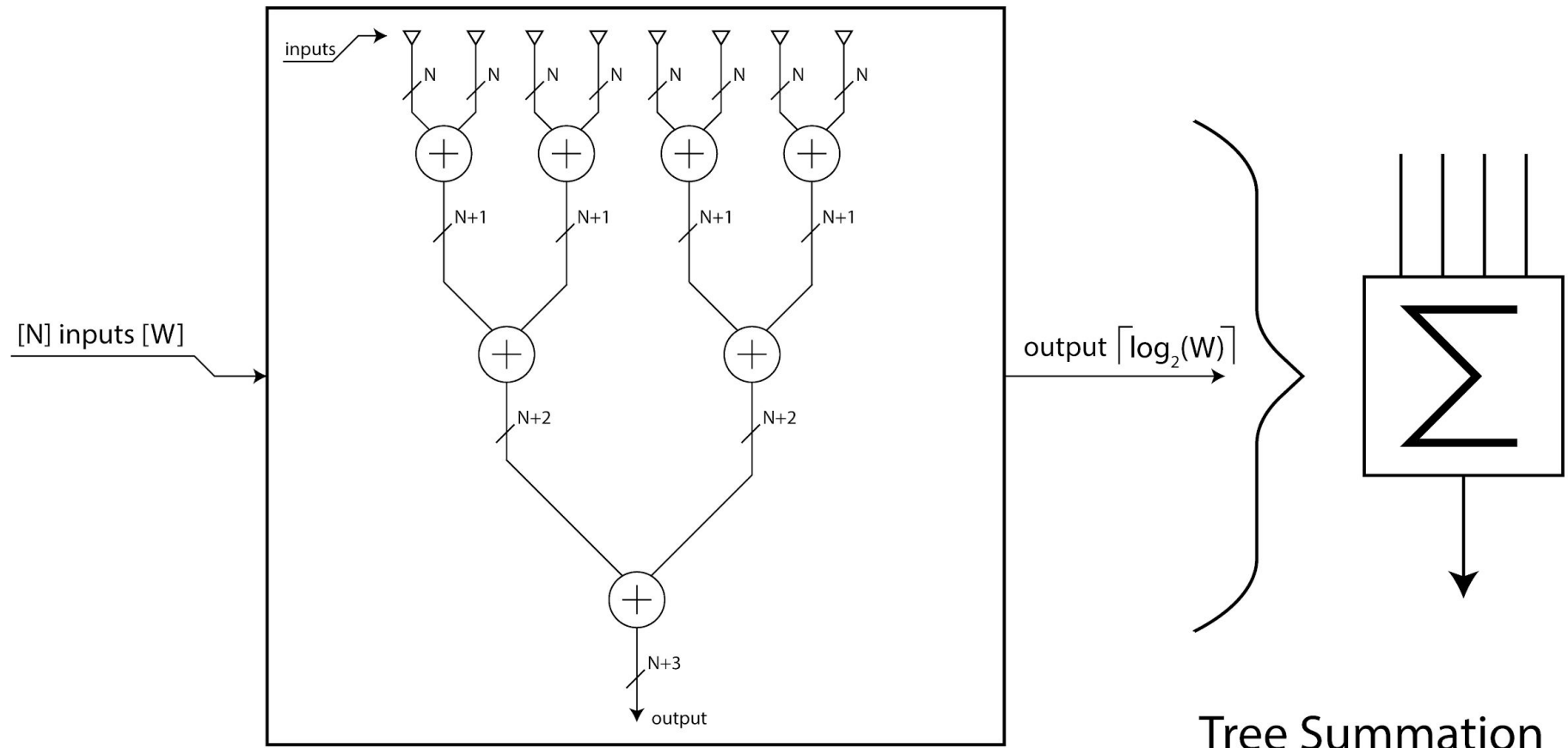
**Risto Miikkulainen** risto@cs.utexas.edu Department of Computer Sciences, The University of Texas at Austin, Austin, TX 78712, USA
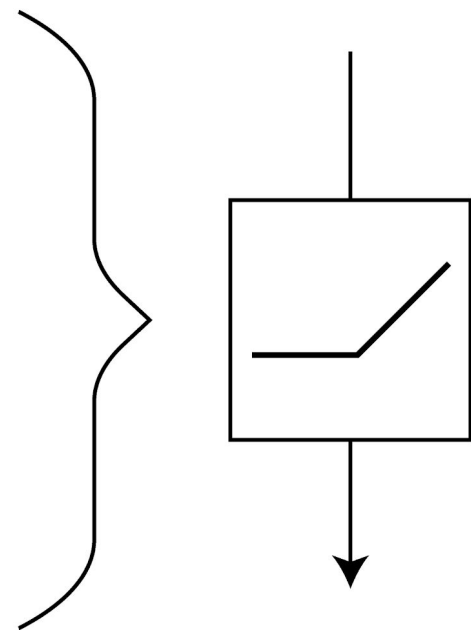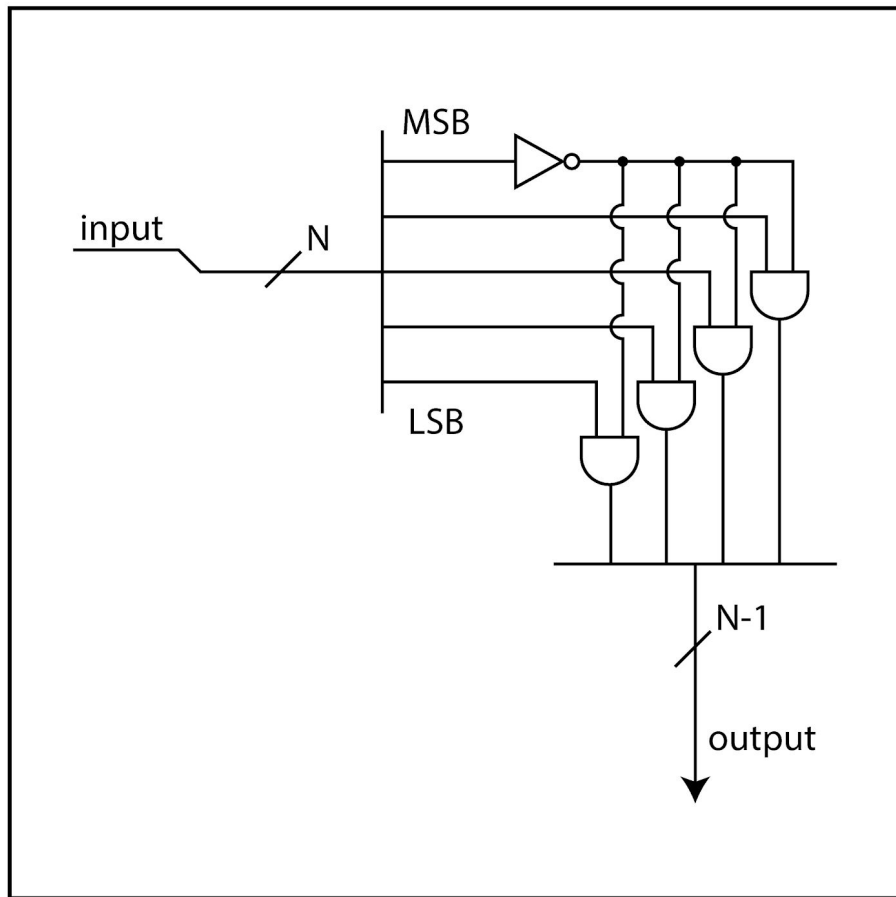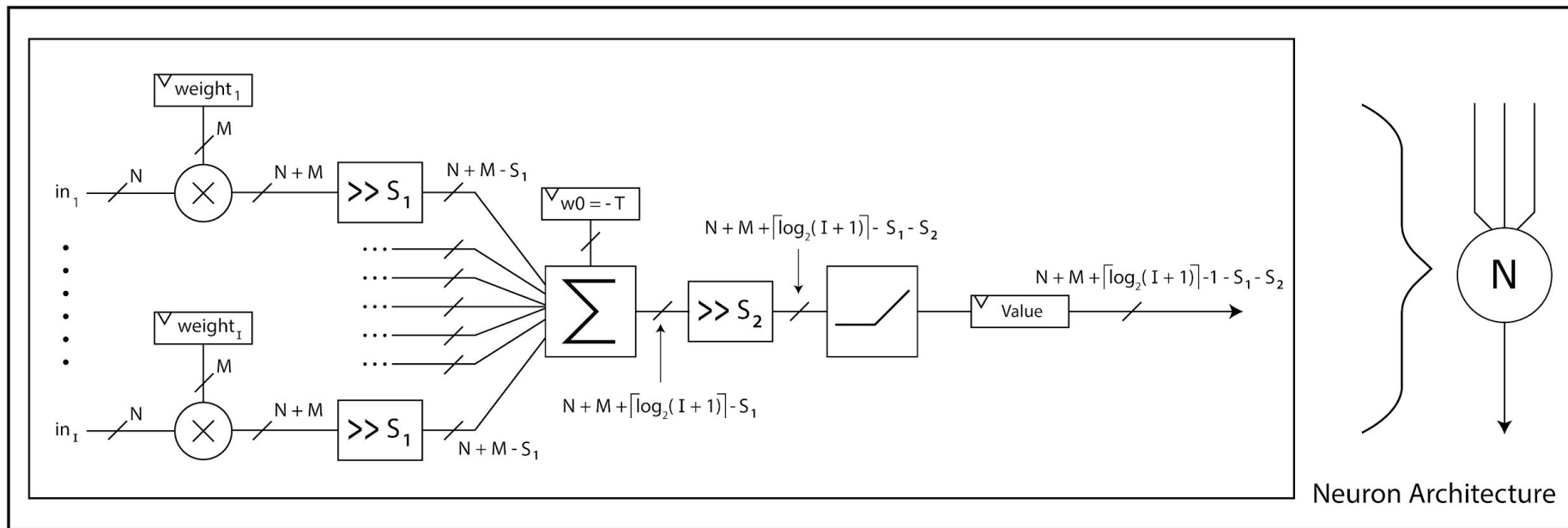
Tree Summation

MSB

input ⟋ N

LSB

N-1

output

Rectified Linear Unit

Neuron Architecture

Input Pool

Communication Controller — data valid — 1

Layer Controller — getFromOutputPool — 1

Communication Controller — data incoming — $I \times N$

Output Pool — new values — $I \times N$

enable

data — $I \times N$

Input Register File

output — $I \times N$

Output Pool

Layer Controller — writeEnable — 1

Layer Controller — dataWidth — $\lceil \log_2 \max(P) \rceil$

Layer Controller — start data address — $\lceil \log_2 I \rceil$

Layer Computation — data value(s) — $N \times T$

Data Register File

data out — $I \times N$

Input Pool

Communication Controller