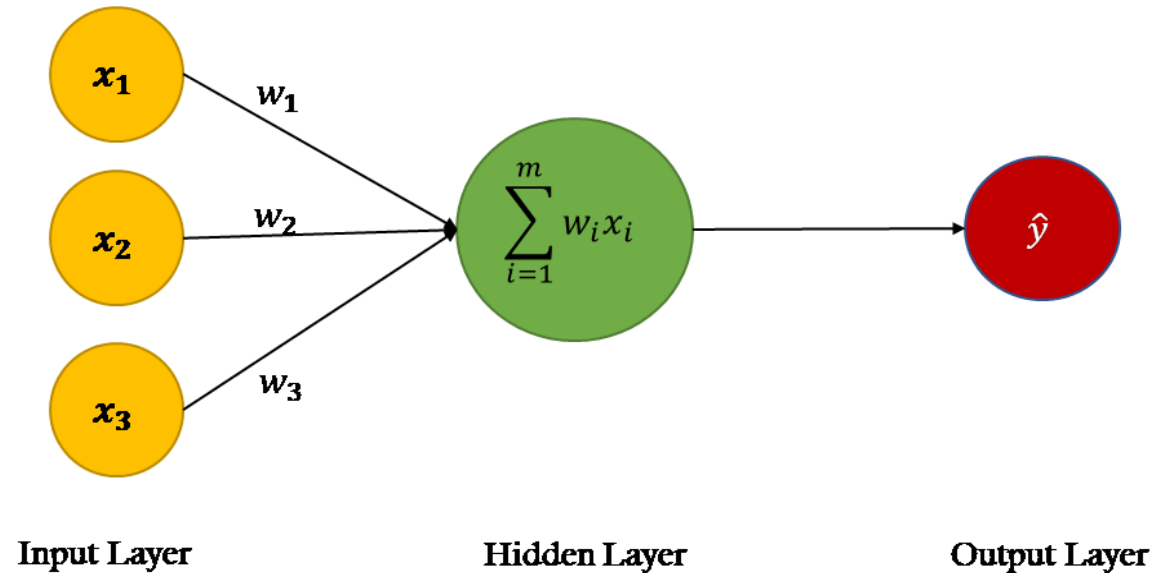# Aplikasi ANN

Erwin Setiawan

# Artificial Neural Network (ANN)

# Artificial Neural Network (ANN) Sederhana

- Contoh ANN sederhana dengan 3 input, 1 hidden node, dan 1 output.



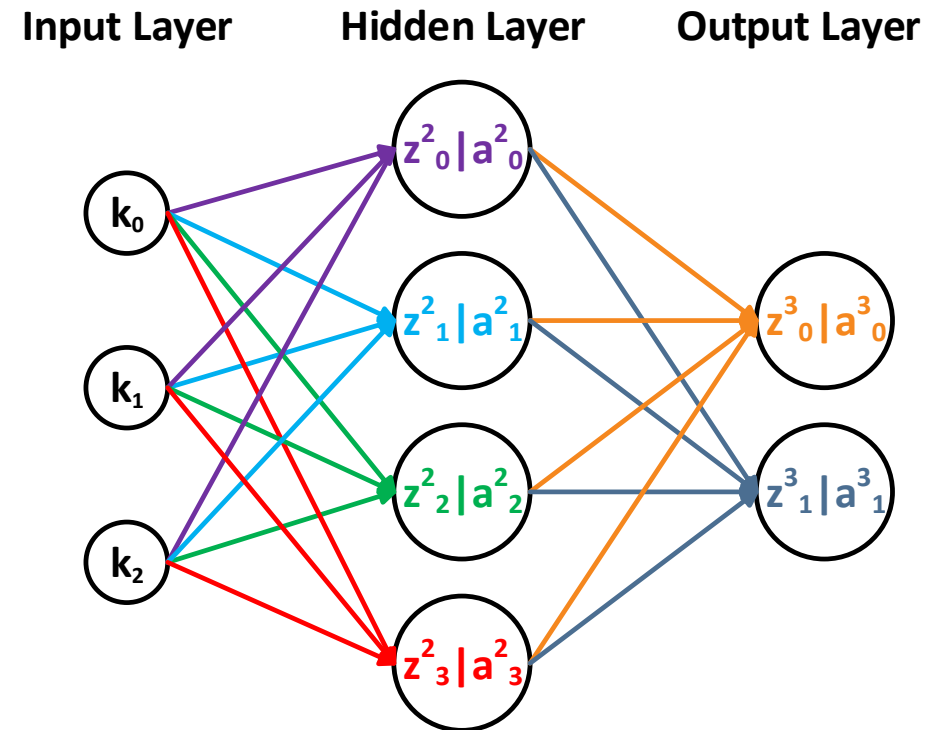$$y = w1*x1 + w2*x2 + w3*x3$$

# ANN dengan Perkalian Matrix

- Menghitung ANN dengan perkalian matrix. Di sini terdapat juga fungsi aktivasi sigmoid ($\sigma$).

$$[k0 \quad k1 \quad k2] \cdot \begin{bmatrix} w00 & w01 & w02 & w03 \\ w10 & w11 & w12 & w13 \\ w20 & w21 & w22 & w23 \end{bmatrix} = [z0 \quad z1 \quad z2 \quad z3]$$

$$\sigma([z0 \quad z1 \quad z2 \quad z3]) = [a0 \quad a1 \quad a2 \quad a3]$$

$$[a0 \quad a1 \quad a2 \quad a3] \cdot \begin{bmatrix} w00 & w01 \\ w10 & w11 \\ w20 & w21 \\ w30 & w31 \end{bmatrix} = [z0 \quad z1]$$

$$\sigma([z0 \quad z1]) = [a0 \quad a1]$$

**Input Layer**    **Hidden Layer**    **Output Layer**
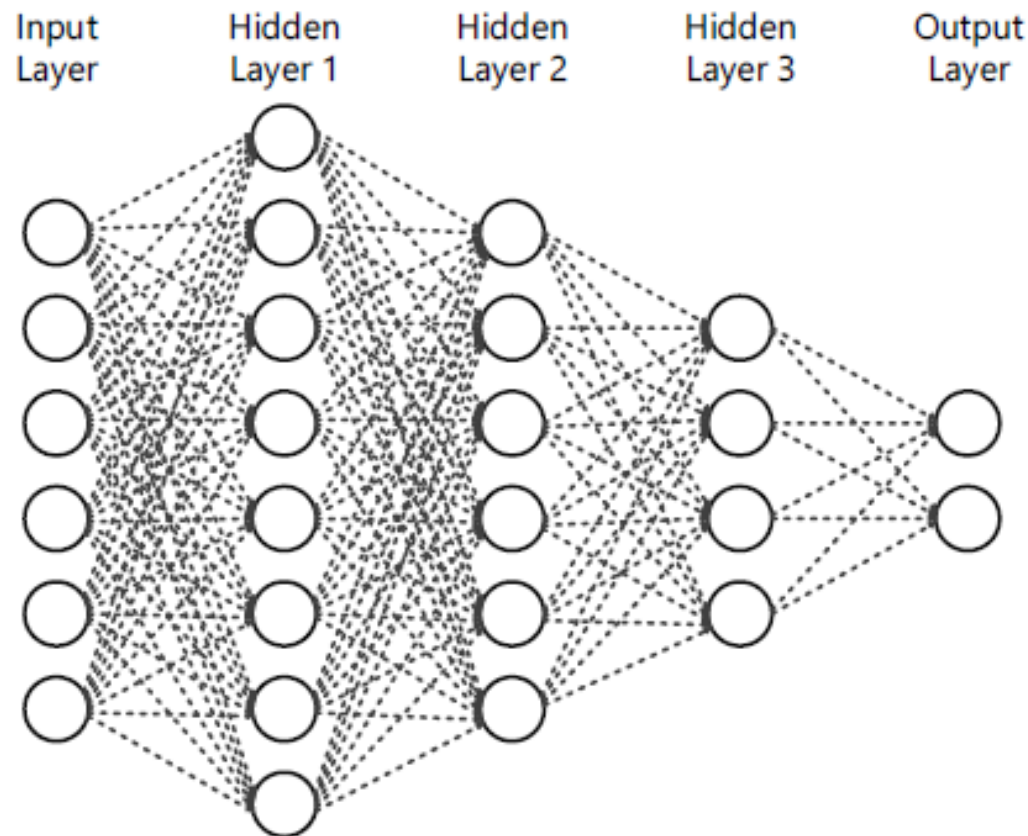
# Deep Neural Network

- ANN yang memiliki hidden layer lebih dari 2 disebut Deep Neural Network (DNN). Contoh model dari paper: Implementation of Systolic Co-processor for Deep Neural Network Inference based on SoC (https://ieeexplore.ieee.org/document/8649920).
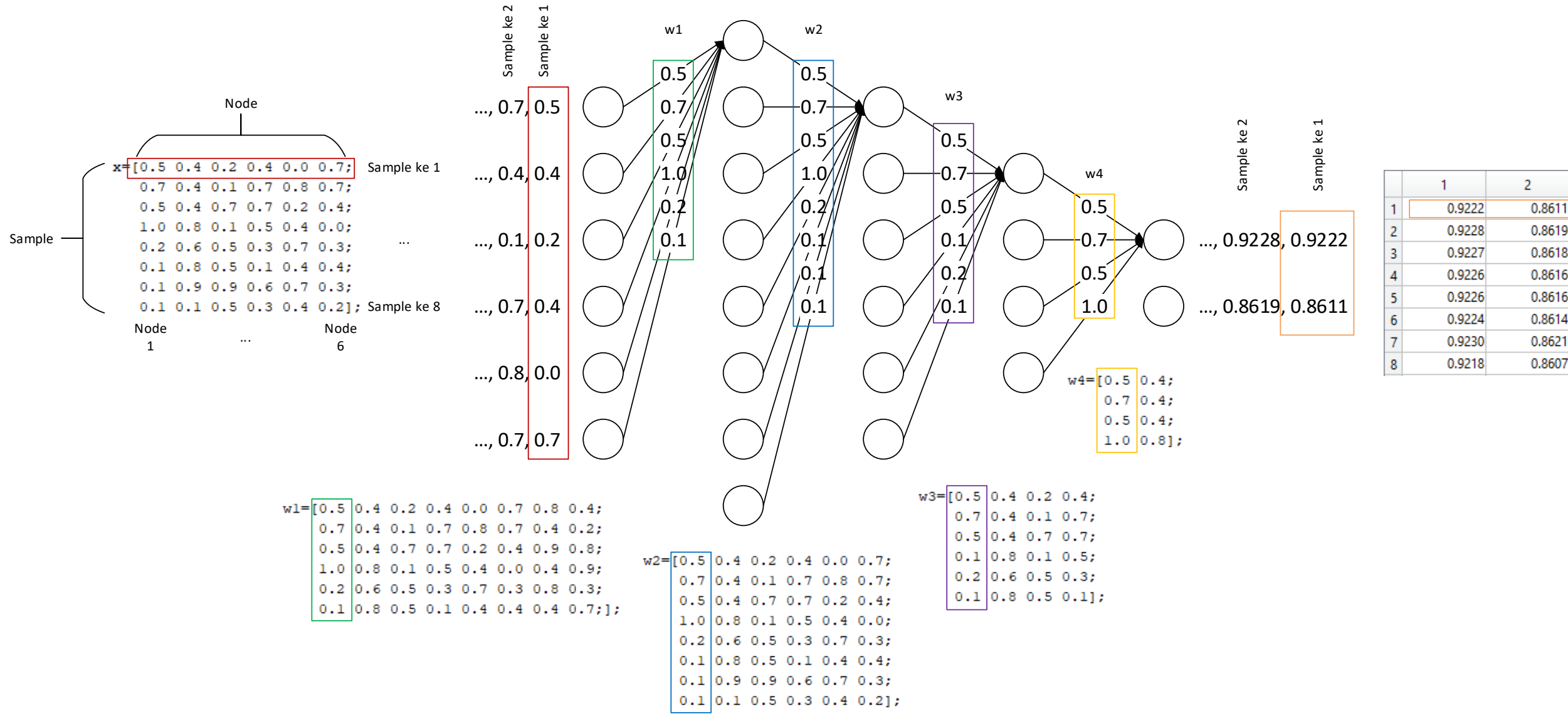
# Model ANN

# Model MATLAB Deep Neural Network

- File model matlab: https://github.com/yohanes-erwin/pemrograman_zynq/blob/main/Aplikasi_ANN/model/model_ann.m
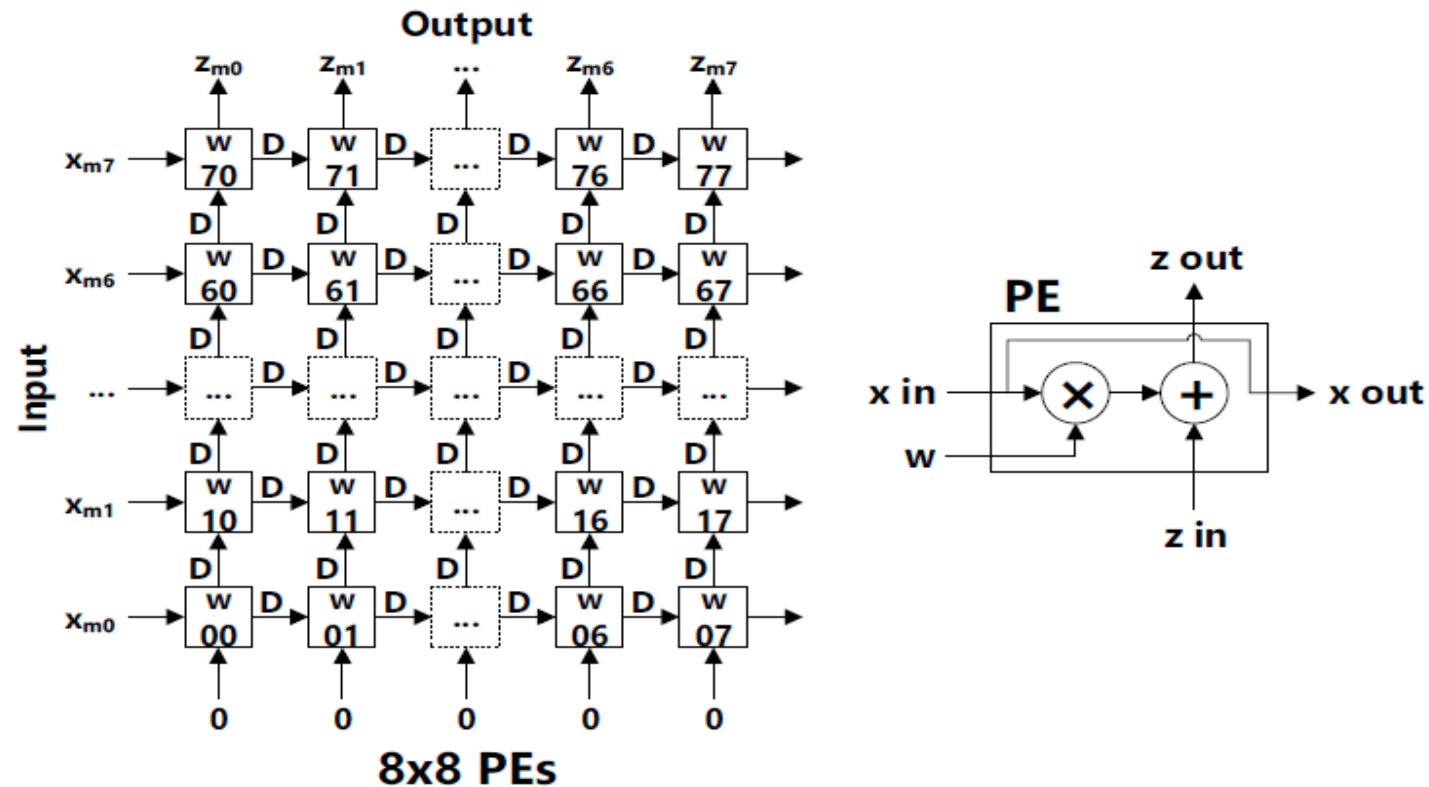
# Ilustrasi Matrix pada Model

# Design Accelerator ANN

# Perkalian Matrix

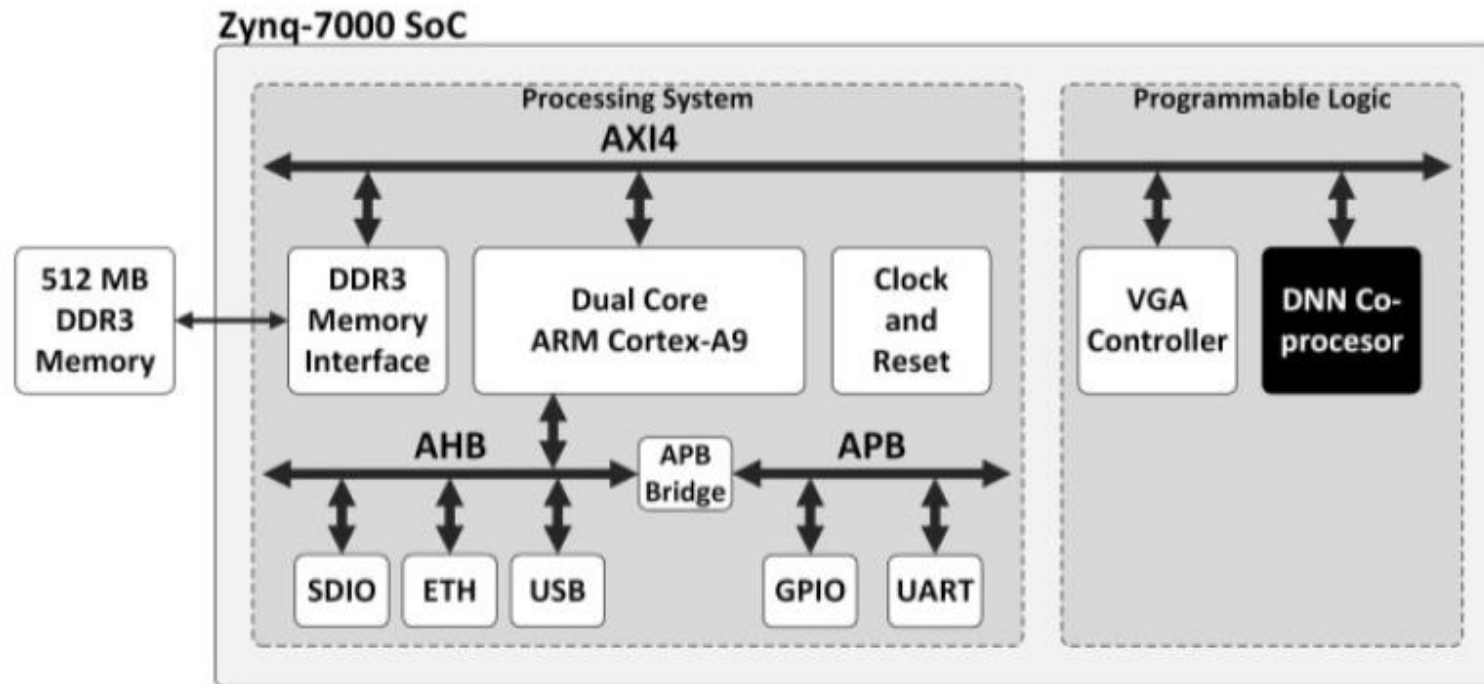- Datapath systolic perkalian matrix untuk ukuran maximal perkalian 8x8.

# Design Accelerator DNN

- Design accelerator (co-processor) DNN dengan datapath systolic matrix 8x8, register untuk data, weight, dan interface AXI4-lite.



**DNN Co-processor**

# Design SoC DNN

- Design SoC secara keseluruhan terdiri dari DNN co-processor dan processing system tempat program C dijalankan untuk mengakses DNN co-processor tersebut.
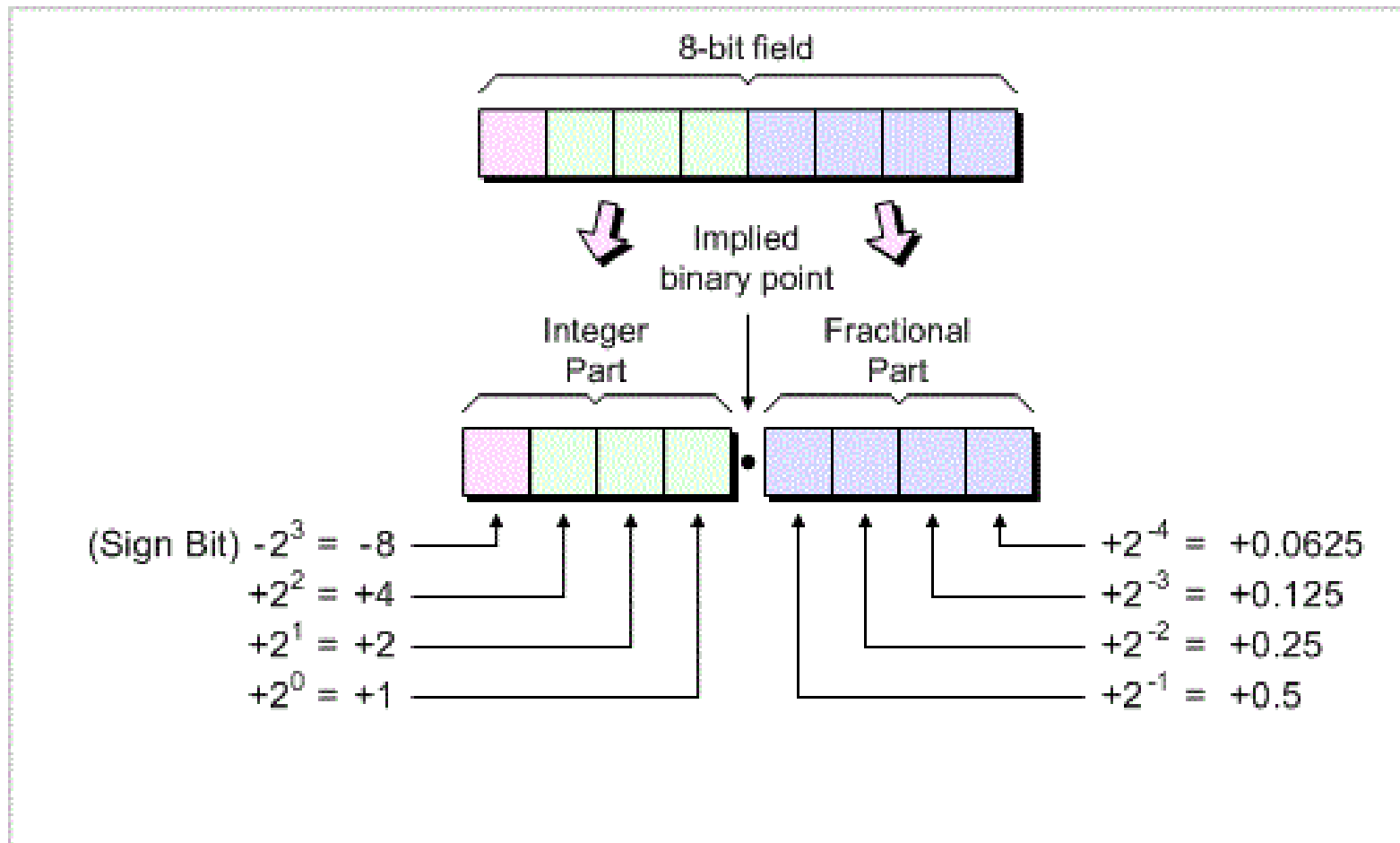
# Kode Verilog Co-processor DNN

- File **module RTL**: https://github.com/yohanes-erwin/pemrograman_zynq/tree/main/Aplikasi_ANN/module

- File **testbench**: https://github.com/yohanes-erwin/pemrograman_zynq/tree/main/Aplikasi_ANN/testbench

- File **program C**: https://github.com/yohanes-erwin/pemrograman_zynq/tree/main/Aplikasi_ANN/program

# Bilangan Fixed Point

# Fixed-Point

# Fixed-Point



- 01010110 = 86 (integer)
- 01010110 = 5.375 (fixed-point 1 sign 3 integer 4 fraction)

$$0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} + 0 \times 2^{-4} = 5.375$$

https://www.allaboutcircuits.com/technical-articles/fixed-point-representation-the-q-format-and-addition-examples/

# Fixed-Point

- Fixed-Point yang digunakan di perkalian marix dan ANN adalah 16-bit dengan
  - 1 sign bit (0: positif, 1: negative)
  - 5 integer bit
  - 10 fraction bit

# Fixed-Point

```
w00 = 16'b00_0010_0001_1001_10;          w00 = 16'b0000010000110011;
```

Tanda _ (underscore) hanya untuk memudahkan pembacaan bit tidak ada pengaruhnya dengan nilainya. Kedua variable tersebut akan bernilai sama. Compiler akan mengabaikan tanda _.

```
w00 = 16'b00_0010_0001_1001_10;
```

Sign (1-bit), Integer (5-bit) | Fraction (10-bit)

•

# Fixed-Point

```
>> q2dec('0000100001100110', 5, 10, 'bin')

ans =

    2.0996

>> dec2q(2.0996, 5, 10, 'bin')

ans =

    '0000100001100101'
```
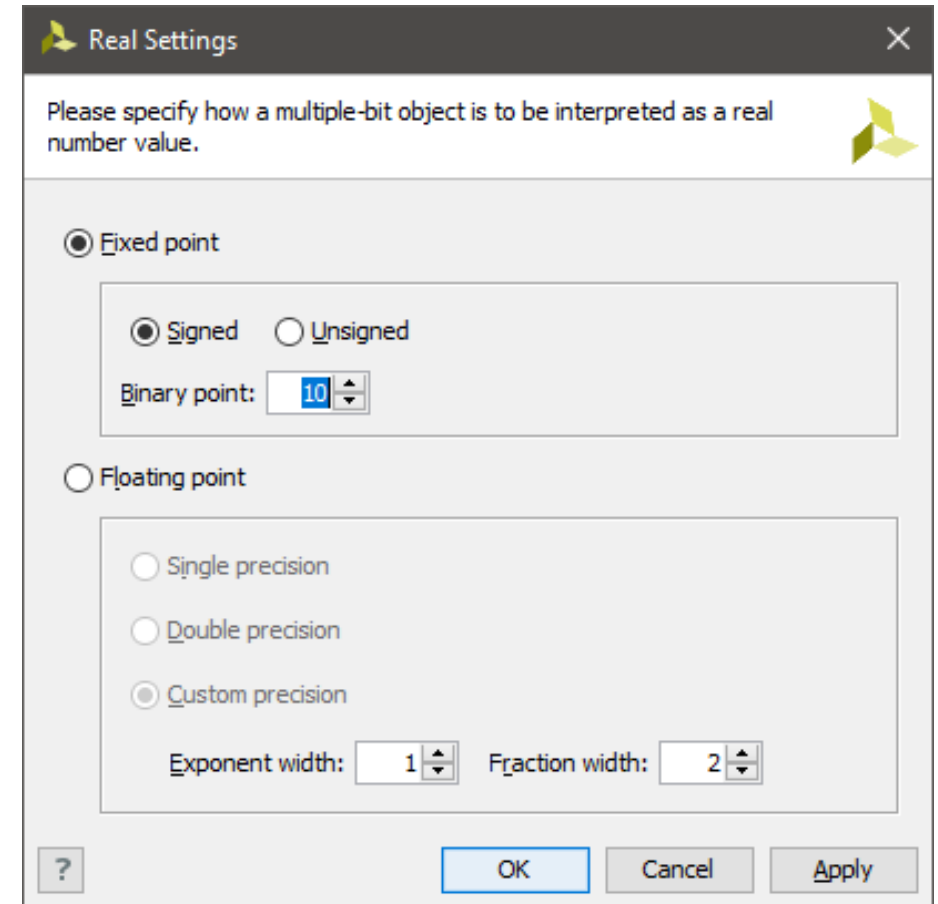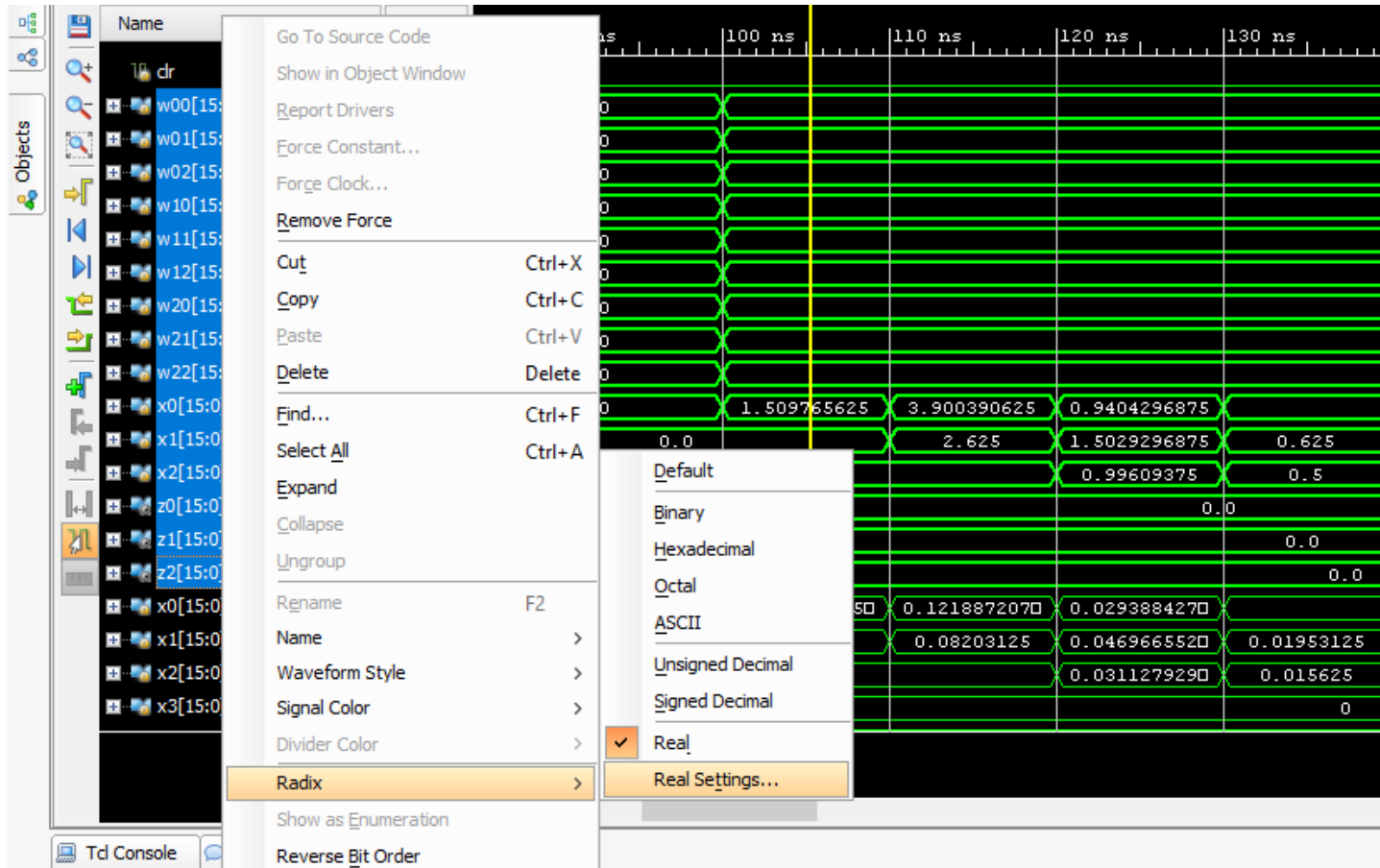
https://www.mathworks.com/matlabcentral/fileexchange/61670-fixed-point-q-format-to-decimal-converter
https://www.mathworks.com/matlabcentral/fileexchange/61669-decimal-to-fixed-point-q-format-converter

# Cara Mengganti Tampilan ke Fixed Point

# Simulasi Perkalian Matrix

# Perkalian Matrix

```
1 -     clear;
2 -     clc;
3 -     x=[1.509 2.625 0.996;
4           3.900 1.502 0.500
5           0.940 0.625 0.519];
6 -     w=[2.099 1.524 0.875;
7           0.527 2.750 1.546;
8           1.945 2.378 0.518];
9 -     z=x*w;
```

```
z =

    6.4880    11.8870     5.8946
    9.9502    11.2631     5.9936
    3.3119     4.3855     2.0576
```

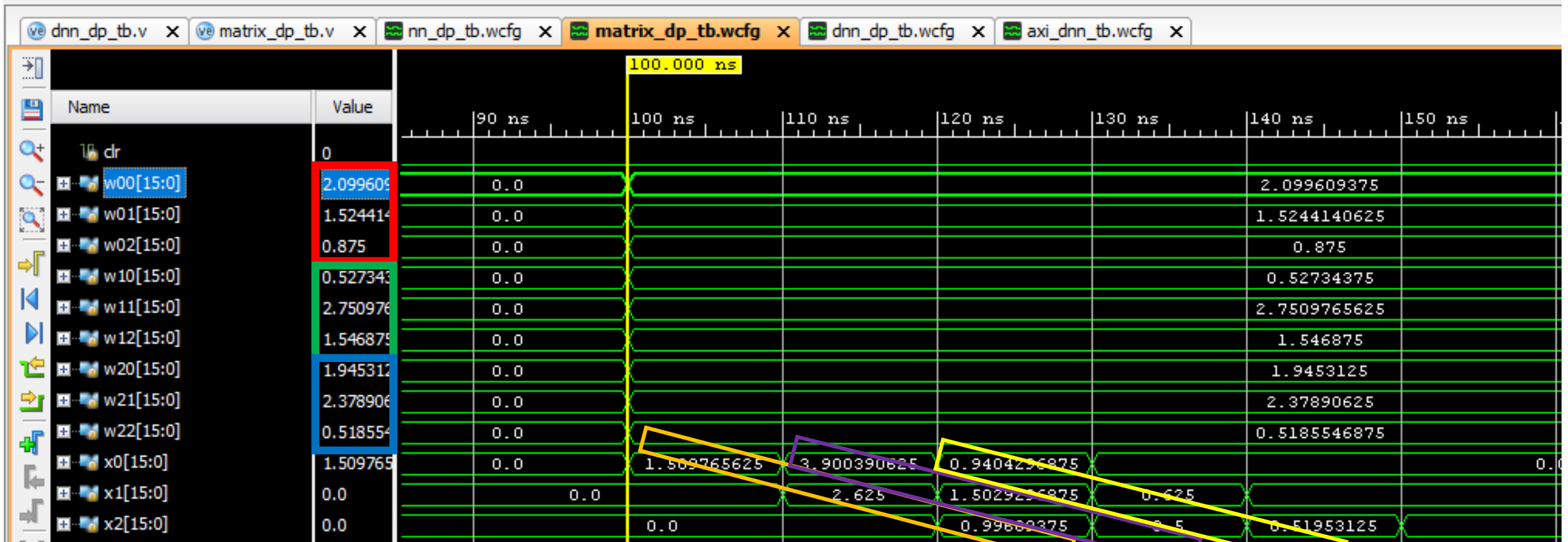# Input

```
1 -    clear;
2 -    clc;
3 -    x=[1.509 2.625 0.996;
4      3.900 1.502 0.500
5      0.940 0.625 0.5191];
6 -    w=[2.099 1.524 0.875;
7      0.527 2.750 1.546;
8      1.945 2.378 0.518];
9 -    z=x*w;
```

# Output

z =

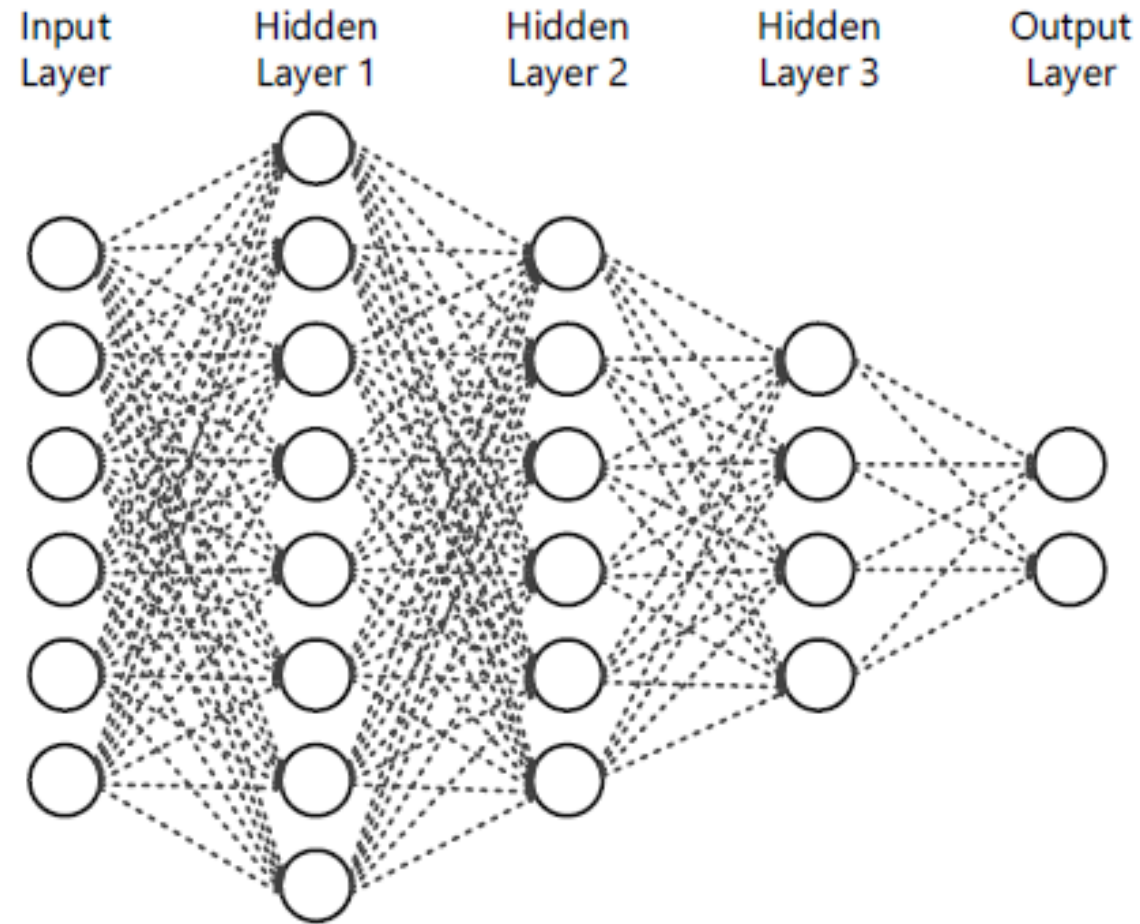| 6.4880 | 11.8870 | 5.8946 |
|--------|---------|--------|
| 9.9502 | 11.2631 | 5.9936 |
| 3.3119 | 4.3855 | 2.0576 |

# Simulasi Proses 1 Layer ANN

# Model ANN



Implementation of Systolic Co-processor for Deep Neural Network Inference based on SoC
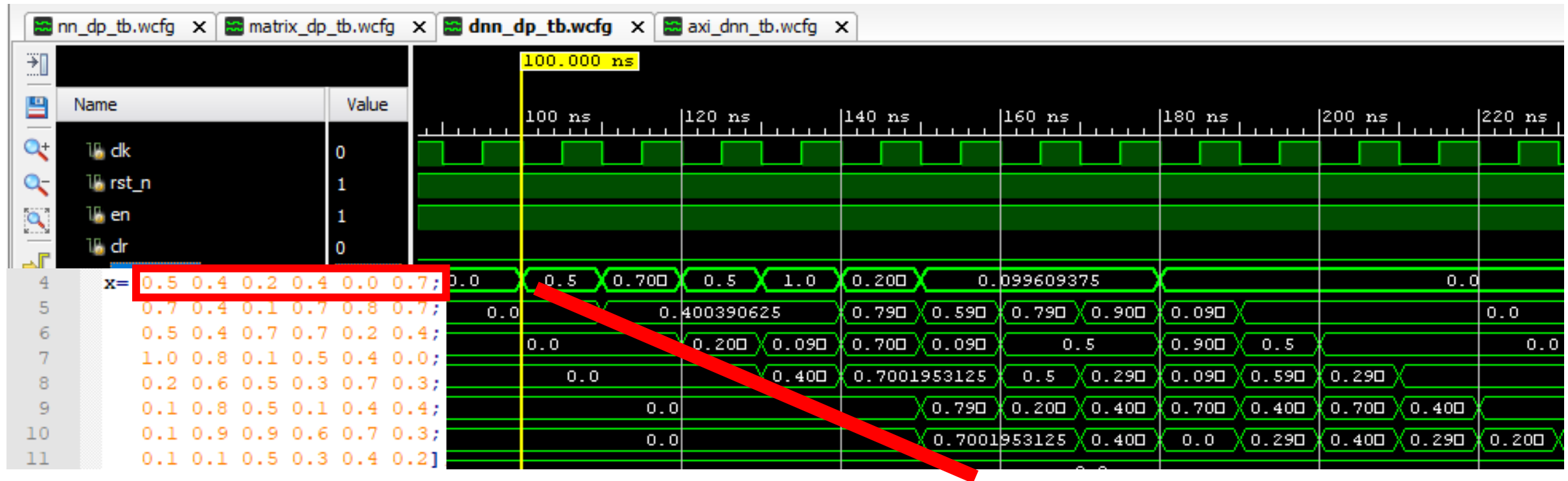
# Model ANN di Matlab

```matlab
clear;
clc;
% Input
x=[0.5 0.4 0.2 0.4 0.0 0.7;
   0.7 0.4 0.1 0.7 0.8 0.7;
   0.5 0.4 0.7 0.7 0.2 0.4;
   1.0 0.8 0.1 0.5 0.4 0.0;
   0.2 0.6 0.5 0.3 0.7 0.3;
   0.1 0.8 0.5 0.1 0.4 0.4;
   0.1 0.9 0.9 0.6 0.7 0.3;
   0.1 0.1 0.5 0.3 0.4 0.2];
% Weight input to hidden 1
w_ih1=[0.5 0.4 0.2 0.4 0.0 0.7 0.8 0.4;
       0.7 0.4 0.1 0.7 0.8 0.7 0.4 0.2;
       0.5 0.4 0.7 0.7 0.2 0.4 0.9 0.8;
       1.0 0.8 0.1 0.5 0.4 0.0 0.4 0.9;
       0.2 0.6 0.5 0.3 0.7 0.3 0.8 0.3;
       0.1 0.8 0.5 0.1 0.4 0.4 0.4 0.7;];
% Weight hidden 1 to hidden 2
w_h1h2=[0.5 0.4 0.2 0.4 0.0 0.7;
        0.7 0.4 0.1 0.7 0.8 0.7;
        0.5 0.4 0.7 0.7 0.2 0.4;
        1.0 0.8 0.1 0.5 0.4 0.0;
        0.2 0.6 0.5 0.3 0.7 0.3;
        0.1 0.8 0.5 0.1 0.4 0.4;
        0.1 0.9 0.9 0.6 0.7 0.3;
        0.1 0.1 0.5 0.3 0.4 0.2];
% Weight hidden 2 to hidden 3
w_h2h3=[0.5 0.4 0.2 0.4;
        0.7 0.4 0.1 0.7;
        0.5 0.4 0.7 0.7;
        1.0 0.8 0.1 0.5;
        0.2 0.6 0.5 0.3;
        0.1 0.8 0.5 0.1];
% Weight hidden 3 to output
w_h3o=[0.5 0.4;
       0.7 0.4;
       0.5 0.4;
       1.0 0.8];
% Output
z1=x*w_ih1;
a1=1./(1+exp(-z1));
z2=a1*w_h1h2;
a2=1./(1+exp(-z2));
z3=a2*w_h2h3;
a3=1./(1+exp(-z3));
z4=a3*w_h3o;
a4=1./(1+exp(-z4));
```

# Simulasi 1 Layer

- dnn_dp_tb.v mensimulasikan satu layer pertama

z1 =

| 1.1000 | 1.3200 | 0.6700 | 0.8900 | 0.8000 | 0.9900 | 1.1800 | 1.2900 |
| 1.6100 | 2.0800 | 1.0700 | 1.2900 | 1.4600 | 1.3300 | 2.0100 | 1.8000 |
| 1.6600 | 1.6400 | 1.0000 | 1.4200 | 1.0400 | 1.1300 | 1.7900 | 1.8100 |
| 1.6900 | 1.4000 | 0.6000 | 1.4000 | 1.1400 | 1.4200 | 1.7300 | 1.2100 |
| 1.2400 | 1.4200 | 0.9800 | 1.2400 | 1.3100 | 1.0900 | 1.6500 | 1.2900 |
| 1.0800 | 1.2000 | 0.8600 | 1.1600 | 1.2200 | 1.1100 | 1.3700 | 1.0900 |
| 1.9000 | 1.9000 | 1.3000 | 1.8400 | 1.7500 | 1.3900 | 2.1700 | 1.9000 |
| 0.7700 | 0.9200 | 0.7100 | 0.7500 | 0.6600 | 0.5400 | 1.0900 | 0.9900 |

$fx$

al =

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 0.7503 | 0.7892 | 0.6615 | 0.7089 | 0.6900 | 0.7291 | 0.7649 | 0.7841 |
| 0.8334 | 0.8889 | 0.7446 | 0.7841 | 0.8115 | 0.7908 | 0.8818 | 0.8581 |
| 0.8402 | 0.8375 | 0.7311 | 0.8053 | 0.7389 | 0.7558 | 0.8569 | 0.8594 |
| 0.8442 | 0.8022 | 0.6457 | 0.8022 | 0.7577 | 0.8053 | 0.8494 | 0.7703 |
| 0.7756 | 0.8053 | 0.7271 | 0.7756 | 0.7875 | 0.7484 | 0.8389 | 0.7841 |
| 0.7465 | 0.7685 | 0.7027 | 0.7613 | 0.7721 | 0.7521 | 0.7974 | 0.7484 |
| 0.8699 | 0.8699 | 0.7858 | 0.8629 | 0.8520 | 0.8006 | 0.8975 | 0.8699 |
| 0.6835 | 0.7150 | 0.6704 | 0.6792 | 0.6593 | 0.6318 | 0.7484 | 0.7291 |

*fx*

# Simulasi Proses Lebih dari 1 Layer

# Block Diagram



DNN Co-processor

- AXI4 Bus
- AXI4-Lite Controller
- Weight Registers 128×16-bit
- I/O Data Registers 48×16-bit
- Control Register
- DNN Controller
- Buffer 64×16-bit
- Weight
- Input
- 8x8 PEs
- Output
- Sigmoid Activation

# Processor Perkalian Matrix

# Jumlah Register Weight



Input Layer · Hidden Layer 1 · Hidden Layer 2 · Hidden Layer 3 · Output Layer

2x4=8

4x6=24

6x8=48

8x6=48

Total weight=
48+48+24+8=128

# Jumlah Register Input



Input Layer    Hidden Layer 1    Hidden Layer 2    Hidden Layer 3    Output Layer

Sample ke-8      Sample ke-1

Total input= 6x8=48

# Jumlah Register Buffer

- Jumlah register buffer berfungsi untuk menyimpan hasil perklaian matrix yang ukurannya menyesuaikan ukuran perkalian matrix.

- Karena ukuran perkalian matrix maximal 8x8, maka jumlah element matrix-nya adalah 8x8=64.
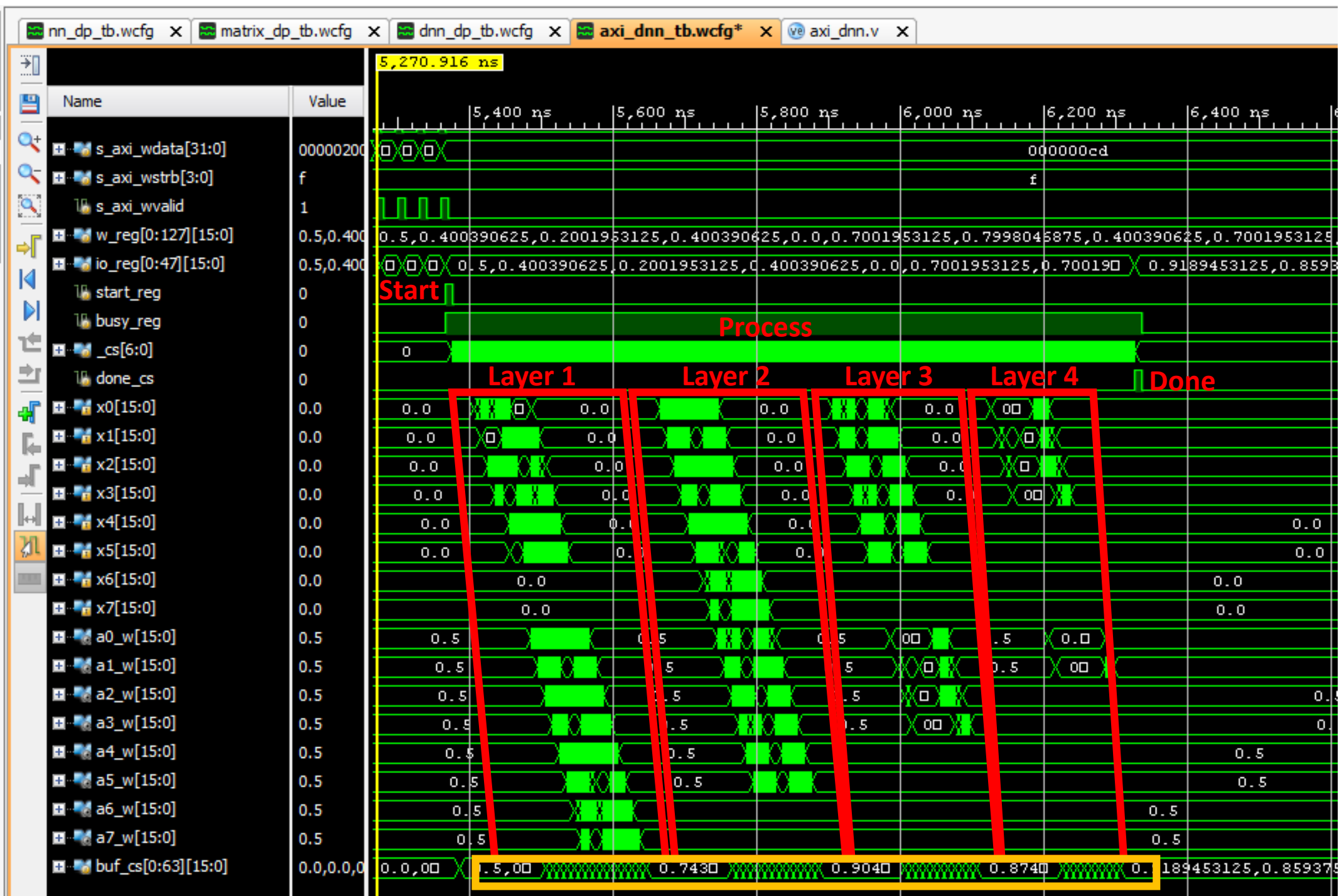
- Sehingga perlu 64 register untuk menyimpan output perkalian matrix.

# State Machine untuk Controller

1. Menunggu sinyal START
2. Load input X dari IO reg ke buf reg
3. Load weight untuk layer n ke weight reg (register systolic)
4. Process perkalian input X dan weight layer n
5. Membaca output perkalian systolic yang sudah dilakukan fungsi sigmoid dan disimpan di buf reg
6. Ulangi ke step 3 sampai jumlah layer 4 terpenuhi
7. Menulis output dari buf reg ke IO reg
8. Generate sinyal DONE
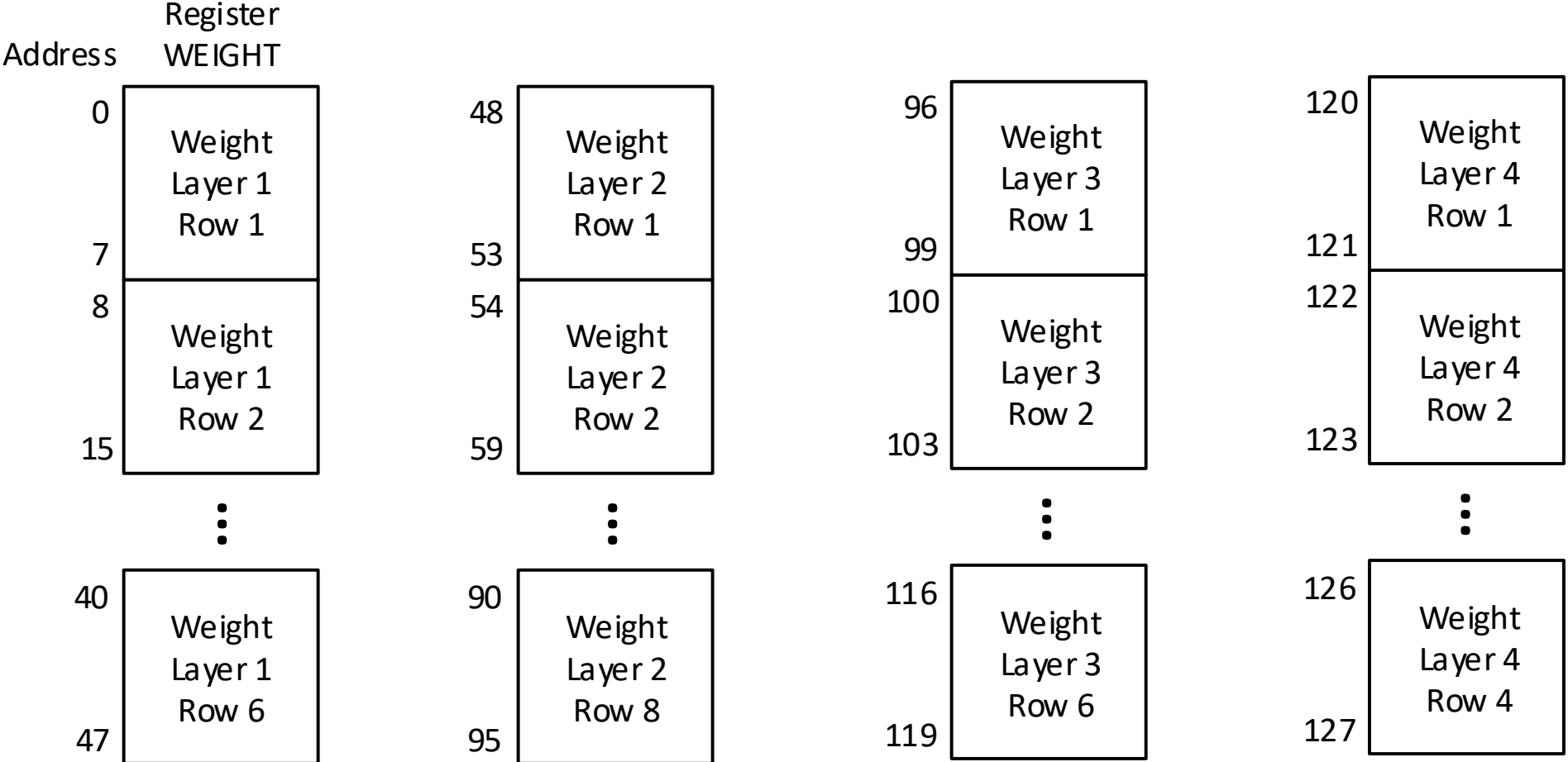
# Address Register Map

# Address Map



1 Reg — Control REG

Bit ke 0 : busy register
Bit ke 31-1: tidak digunakan

128 Reg — Weight REG

48 Reg — IO REG

Start core secara otomatis berjalan ketika
menulis data terakhir ke register IO alamat 48

# Address Register Weight

Address

Register WEIGHT

| 0 | Weight Layer 1 Row 1 |
| 7 | |
| 8 | Weight Layer 1 Row 2 |
| 15 | |

⋮

| 40 | Weight Layer 1 Row 6 |
| 47 | |

| 48 | Weight Layer 2 Row 1 |
| 53 | |
| 54 | Weight Layer 2 Row 2 |
| 59 | |

⋮

| 90 | Weight Layer 2 Row 8 |
| 95 | |

| 96 | Weight Layer 3 Row 1 |
| 99 | |
| 100 | Weight Layer 3 Row 2 |
| 103 | |

⋮

| 116 | Weight Layer 3 Row 6 |
| 119 | |

| 120 | Weight Layer 4 Row 1 |
| 121 | |
| 122 | Weight Layer 4 Row 2 |
| 123 | |

⋮

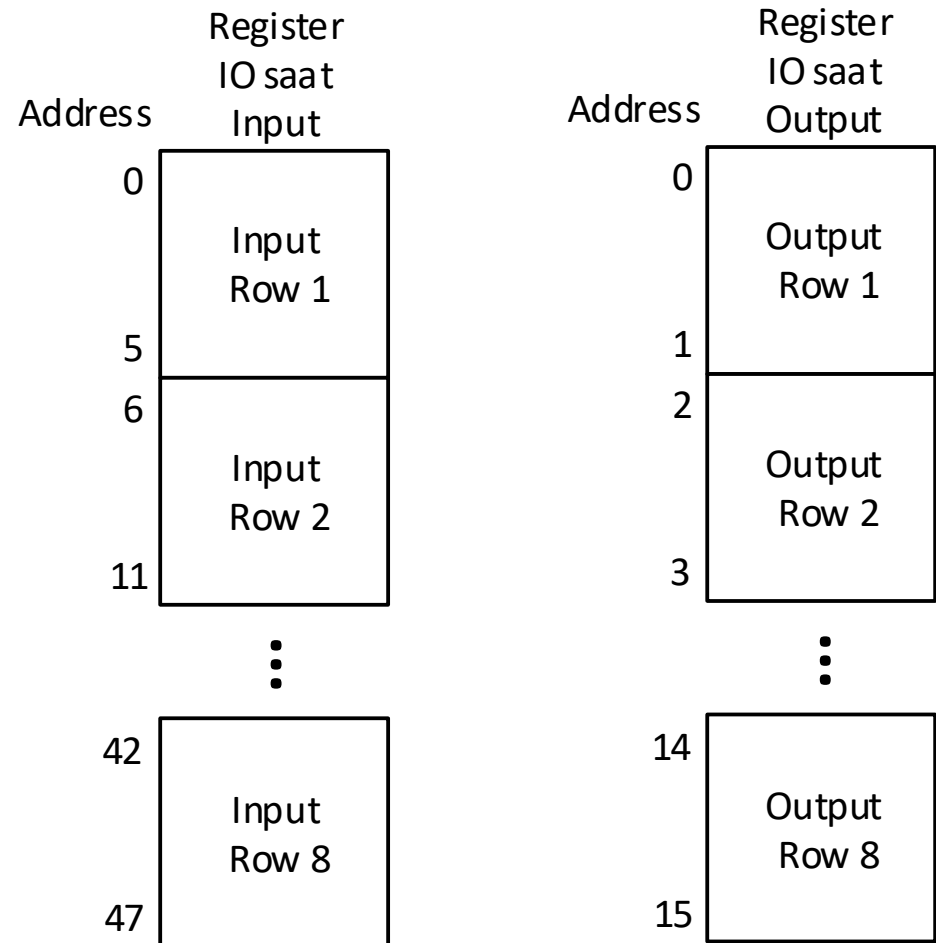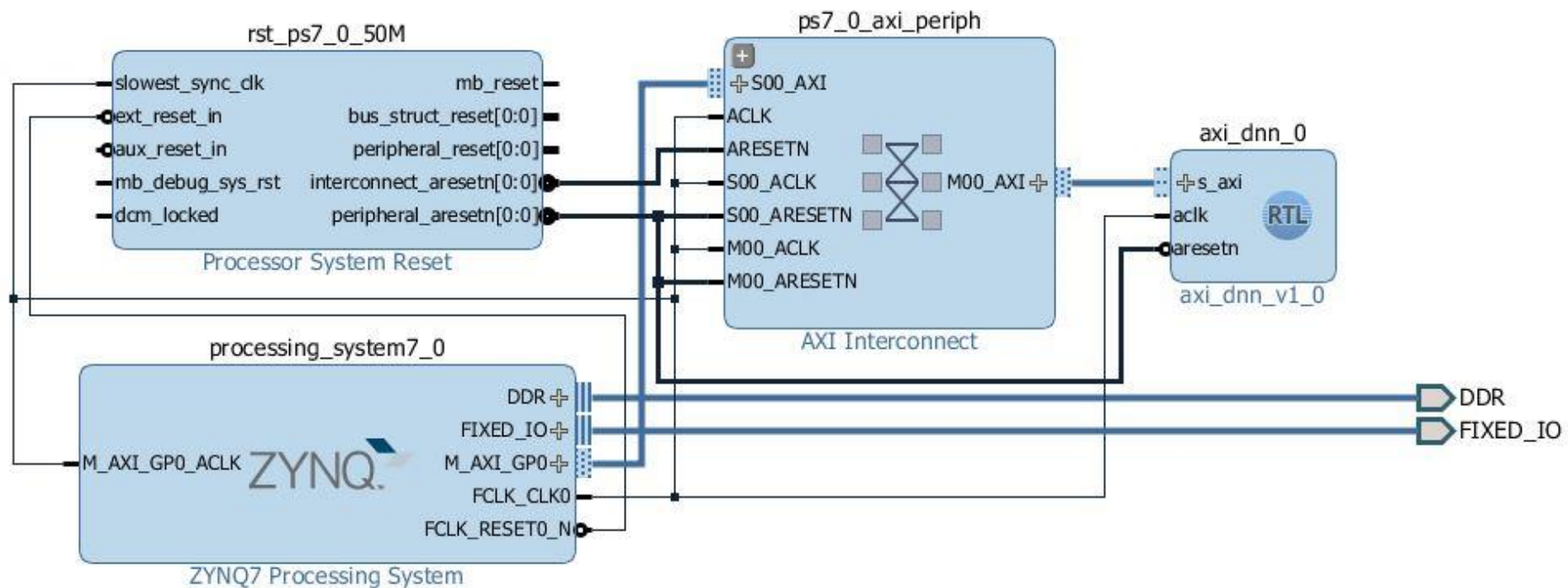| 126 | Weight Layer 4 Row 4 |
| 127 | |

# Address Register IO

# Implementasi di Board ZYBO

# Block Design

- Tambahkan ZYNQ PS dan aktifkan UART 1
- Tambahkan module axi_dnn

# Output SDK Terminal



Connected to: Serial ( COM10, 115200, 0, 8 )

DNN start =====
w1: 512, 410, 205, 410, 0, 717, 819, 410, 717, 410, 102, 717, 819, 717, 410, 205, 512, 410, 717, 717, 205, 410, 922, 819, 1024, 819, 102, 512, 410, 0, 410, 922, 205, 614, 512, 307, 717, 307, 819, 307, 102, 819, 512, 102, 417, 417, 417, 717,
w2: 512, 410, 205, 410, 0, 717, 717, 410, 102, 717, 819, 717, 512, 410, 717, 717, 205, 410, 1024, 819, 102, 512, 410, 0, 205, 614, 512, 307, 717, 307, 102, 819, 512, 102, 410, 410, 102, 922, 922, 614, 717, 307, 102, 102, 512, 307, 410, 205,
w3: 512, 410, 205, 410, 717, 410, 102, 717, 512, 410, 717, 717, 102, 819, 102, 512, 205, 614, 512, 307, 102, 819, 512, 102,
w4: 512, 410, 717, 410, 512, 410, 1024, 819,
In: 512, 410, 205, 410, 0, 717, 717, 410, 102, 717, 819, 717, 512, 410, 717, 717, 205, 410, 1024, 819, 102, 512, 410, 0, 205, 614, 512, 307, 717, 307, 102, 819, 512, 102, 410, 410, 102, 922, 922, 614, 717, 307, 102, 102, 512, 307, 410, 205,
Busy flag: 0
Out: 941, 880, 941, 880, 941, 880, 941, 880, 941, 880, 941, 880, 941, 880, 941, 872,
Out (fixed-point): 0.918945, 0.859375, 0.918945, 0.859375, 0.918945, 0.859375, 0.918945, 0.859375, 0.918945, 0.859375, 0.918945, 0.859375, 0.918945, 0.859375, 0.918945, 0.851562,
DNN end =====

Send    Clear

# Thank You