# MNIST - Accelerator :

→ Design Flow

```
┌──────────────┐        ╭─────────────╮
│ Python       │────────│  NN Arch    │
│ Model        │        ╰─────────────╯
└──────────────┘
       │
       │  export weights, bias
       ▼
┌──────────────────────┐
│ Datapath /           │◄───────────────┐
│ Pipeline  RTL        │                │
└──────────────────────┘                │
       │                                │
       ▼                                │
┌──────────────────────┐                │
│ Sim + Verification    │◄──────────────┤
└──────────────────────┘                │
       │                                │
       ▼                                │
┌──────────────────────┐                │
│ Vivado  Synth        │                │
└──────────────────────┘                │
       │────────────────────────────────┘
       ▼
┌────────────────────────────────────┐
│ Wrappers + Implementation          │
└────────────────────────────────────┘
       │
       ▼
┌──────────────────────┐
│ Validation           │
└──────────────────────┘
       │
       ▼
┌──────────────────────┐
│ Documentation        │
└──────────────────────┘
```

╭─────────────────────────────────────╮
│  Zynq - 7000      Stimulus          │
│      injection                      │
╰─────────────────────────────────────╯

# MNIST Neural Network Design

Input : 28 x 28 , 8-bit grayscale images flattened to a 784 vector.

Output : Size 10 vector representing strength for each digit.

Architecture : $784 \rightarrow 500 \rightarrow 10$
    Linear (784, 500)
    ReLU (500)
    Linear (500, 10)

Linear Layers : $y = x A^T + b$

    Layer #1 : $A = (500, 784)$
              $b = (500)$

    Layer #2 : $A = (500, 10)$
              $b = (10)$
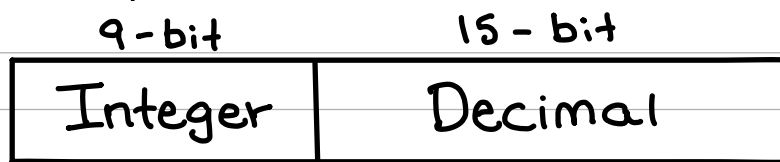
Neurons : $y = xw + b$

    Layer #1 : $x = (784)$
              $w = (784)$
              $b = (1)$
              $y = (1)$

    Layer #2 : $x = (500)$
              $w = (500)$
              $b = (1)$
              $y = (1)$

# RTL Design:

## Key Problems:
→ Data Representation (Bits? Fixed?)
→ Data Operations (Add, Sub, Mult)
→ Functions
- ReLU $(\max(0, z))$
  + Neuron / Layer Design
- Linear $(y = x A^T + b)$
  + Neuron / Layer Design
→ Pipelining / Datapath

## 24 - Bit Fixed Point:

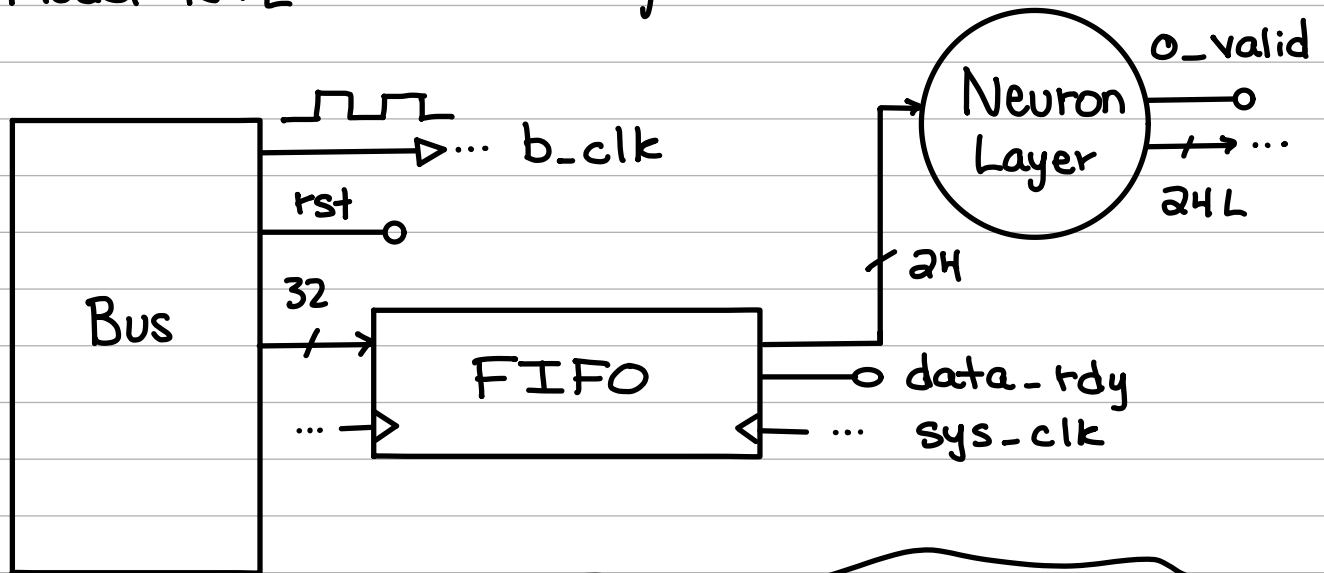| 9-bit | 15 - bit |
|---------|----------|
| Integer | Decimal |

- Add, Sub are normal
- Mult must >> by # decimal (30)

## ReLU Cell RTL:



$z$ ──── 0
'0 ──── 1 ──── ReLU($z$)

$z[23]$

# Model RTL:   L = # Layers

Linear #1



Bus

b_clk

rst

32 → FIFO

Neuron Layer — o_valid

24L

24 → data_rdy / sys_clk

---

Linear #1

Neuron Layer — o_valid

24L

ReLU Nonlinear Layer — o_valid

24L

L 24-bit Register File

i_valid → L cycle ctrl

24

24

Neuron Layer — o_valid

24·10

Argmax 0...9 — final_out / 4

---

Linear #2

b_clk

rst

Bus

final_out / 4 → FIFO

sys_clk ...

32 → data_rdy / b_clk

# Linear Layer RTL:



# Neuron RTL: