

CNN FPGA System Architecture Review

System Overview

- Multi-model CNN implementation (10 models)
- 16x12 image input processing
- Target frequency: 73MHz (non-pipelined) / 51MHz (pipelined)
- 16 parallel compute engines

Key Components & Architecture

1. Processing Core:

- 16 Parallel Compute Engines
- Each engine: 7 DSP blocks (3x 9x9 multipliers, 4x 19x8 multipliers) total of 112 DSP to run one column at once.
- The result of the MultiMultiplierEngine will be ready by the next CLK
- ReLU activation implementation
- Fully connected layer processing

2. Memory Architecture:

- CONV_W_MEM: Convolution weights $[15:0] \times N$
- CONV_B_MEM: Convolution bias $[15:0] \times N$
- WFC_MEM: FC weights $[127:0] \times 3 \times N$
- BFC_MEM: FC bias $[44:0] \times N$
- Relu_out_FIFO $[18:0] \times 128 \times 3 \times N$
- P1xel_MEM int8 $[15:0] \times \text{MAX_W}$

3. Data Management:

- Triple buffer design for continuous processing of each line
- Optimized memory controller
- Efficient address management
- Pipeline-optimized data flow
- Correctly only the 2 line buffers and the FIFO need EN

Critical Parameters

- T_RAM: RAM request load time
- N: number of models
- T_find_winner: Winner calculation time
- MAX_W: the maximum length that an image can be we from correct testing the longest image was 350

Expected Performance for testing MultiMultiplierEngine

- Processing speed: Up to 73MHz
- Parallel model execution
- Continuous data flow capability

Implementation Considerations

1. Resource Requirements:

- DSP blocks: 112 (7 × 16)
- Memory banks: 6 main types

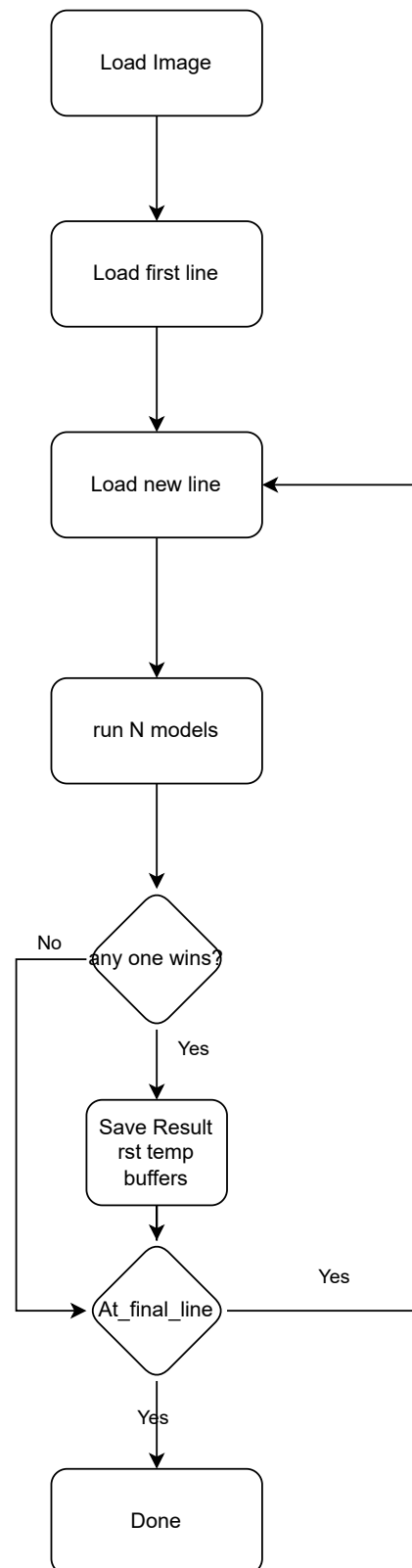
2. Critical Paths:

- Memory access timing
- Compute engine pipeline
- Winner determination logic

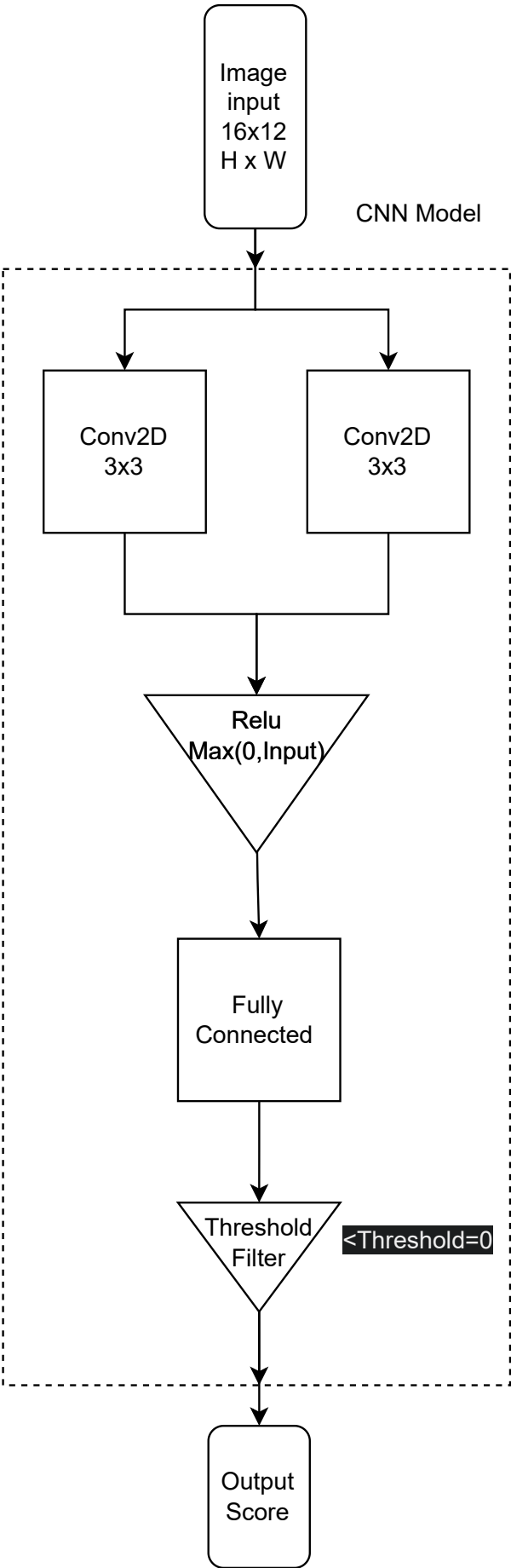
3. Potential Challenges:

- Memory bandwidth optimization
- Pipeline synchronization
- Resource balancing

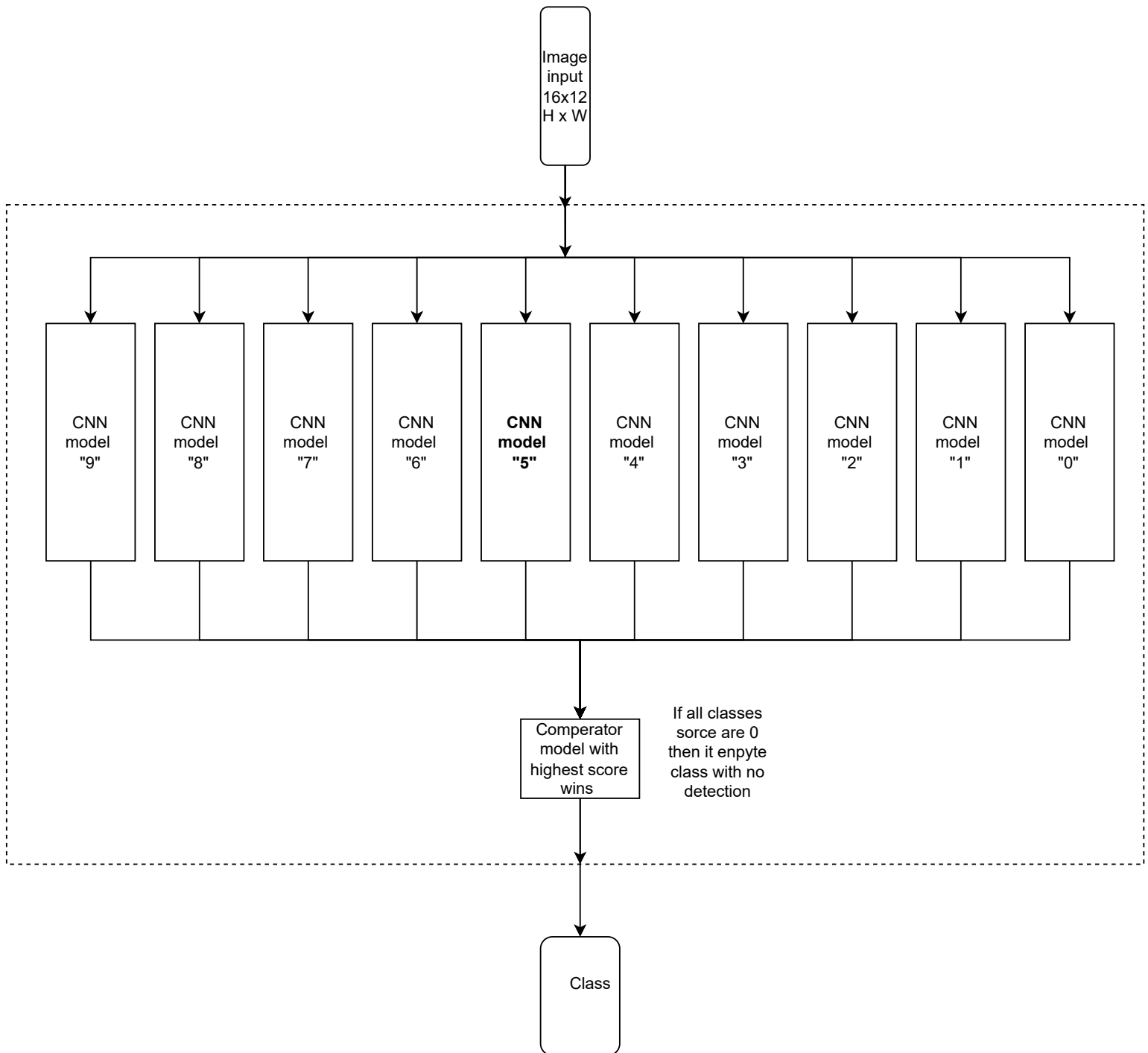
System WorkFlow



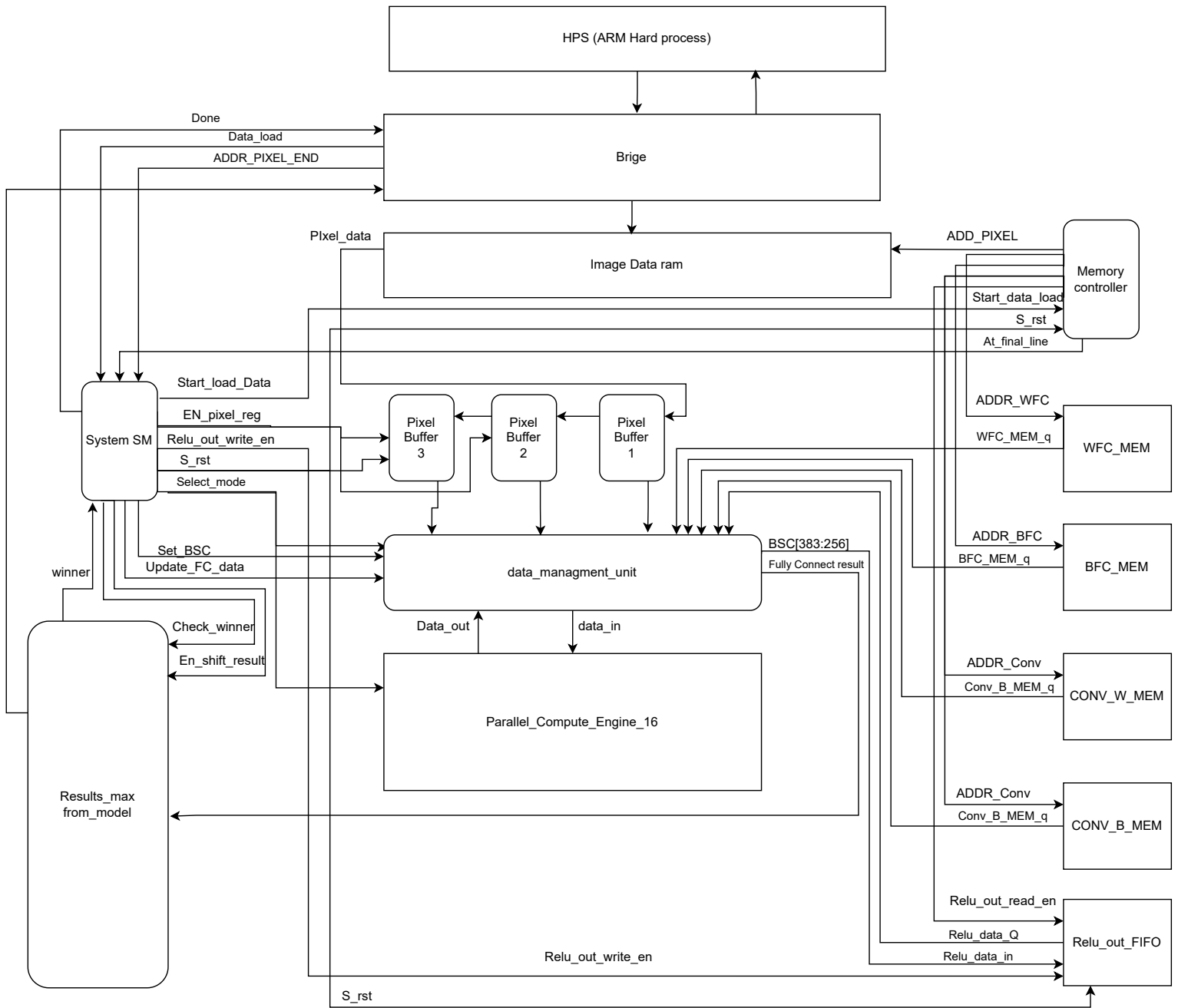
Single CNN network



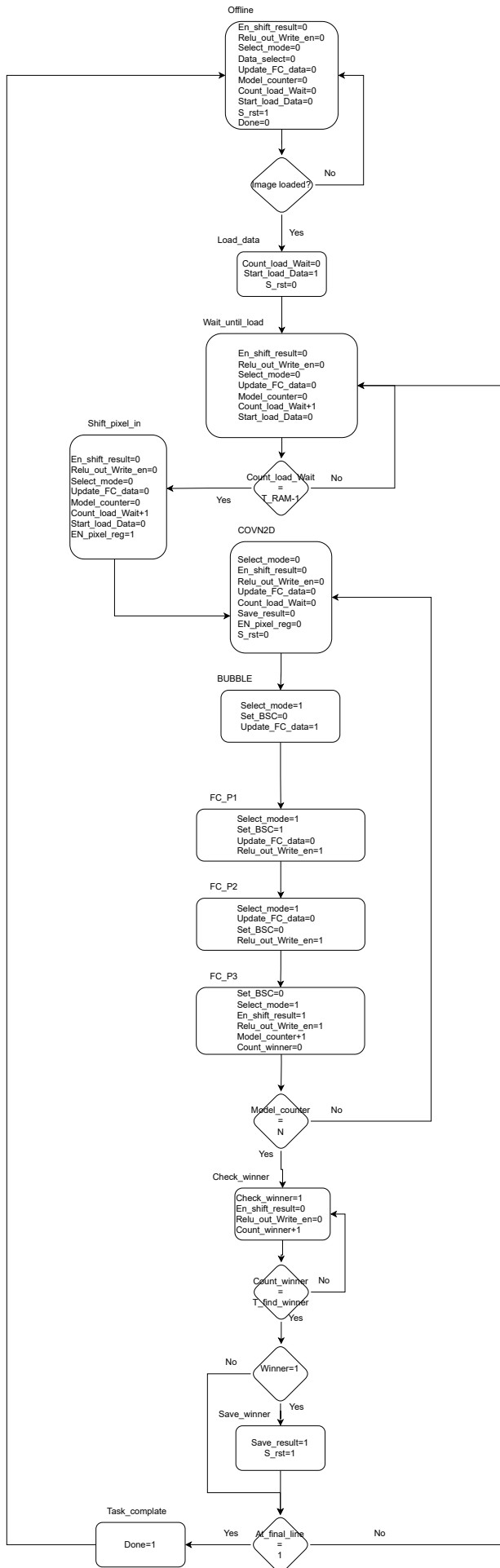
Full system CNN workflow



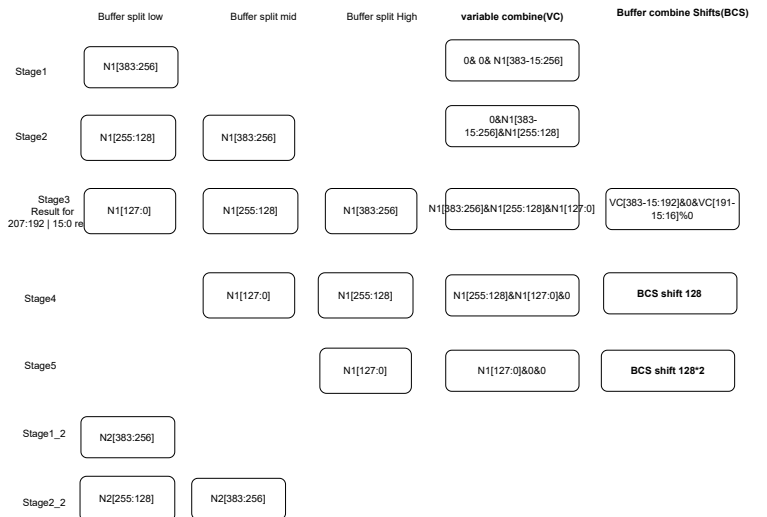
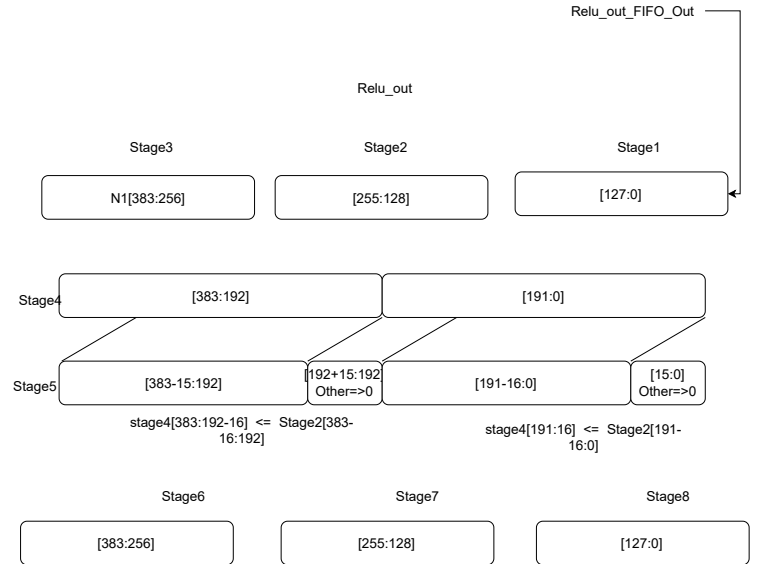
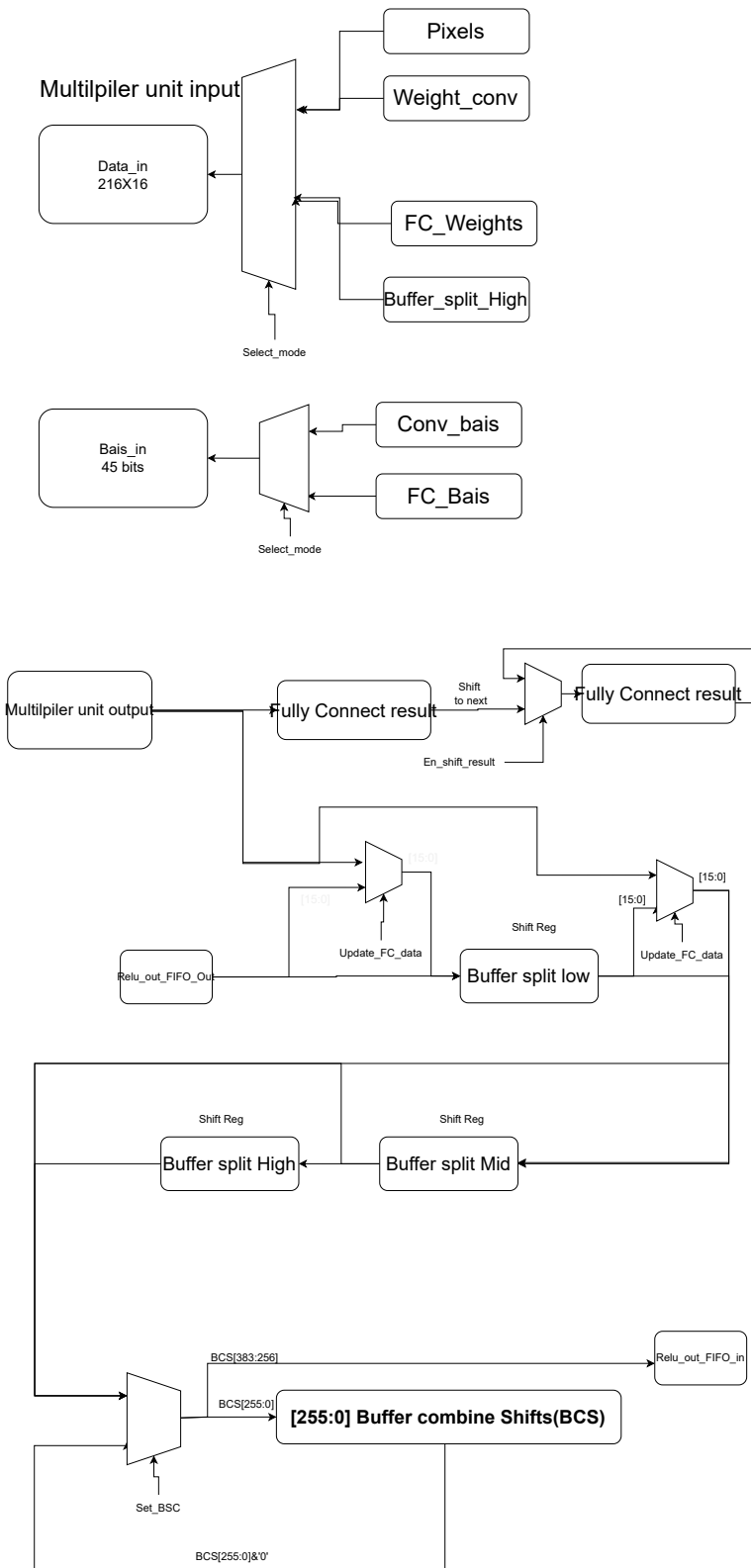
System planed Architecture



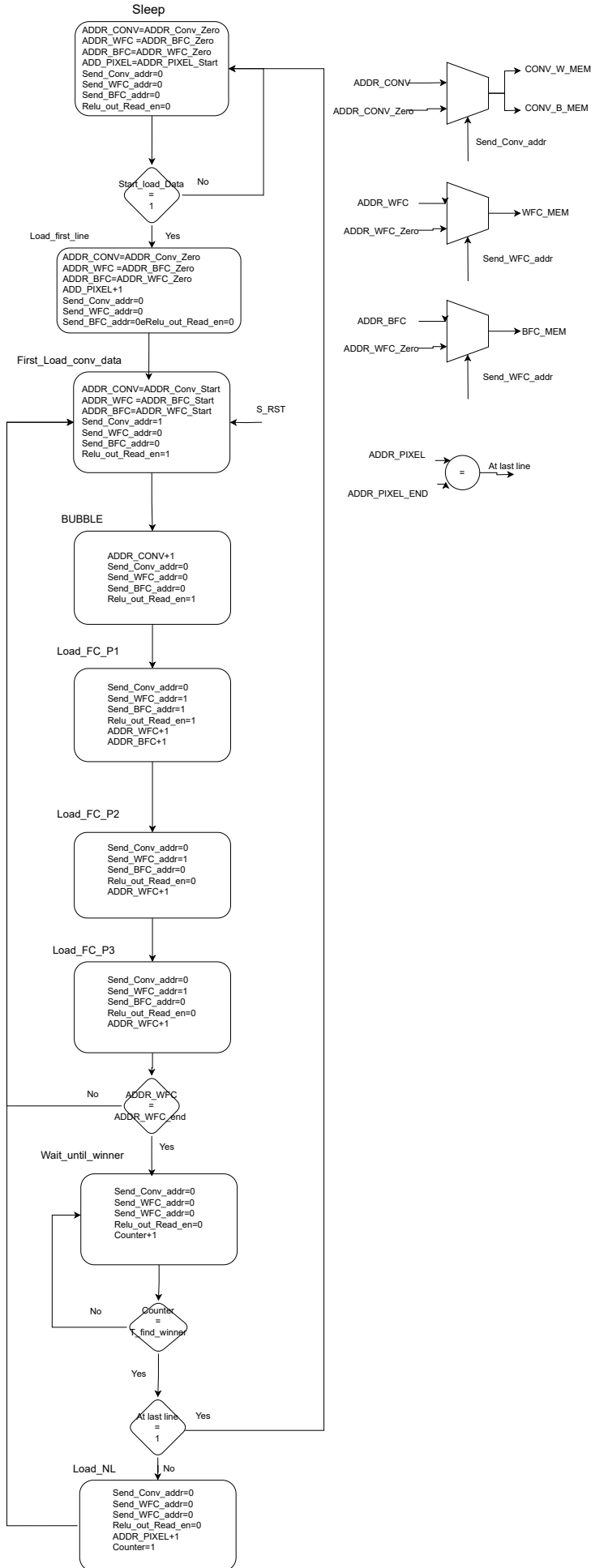
Main system SM



Input buffer load from
Memory controller and manages by Main system SM

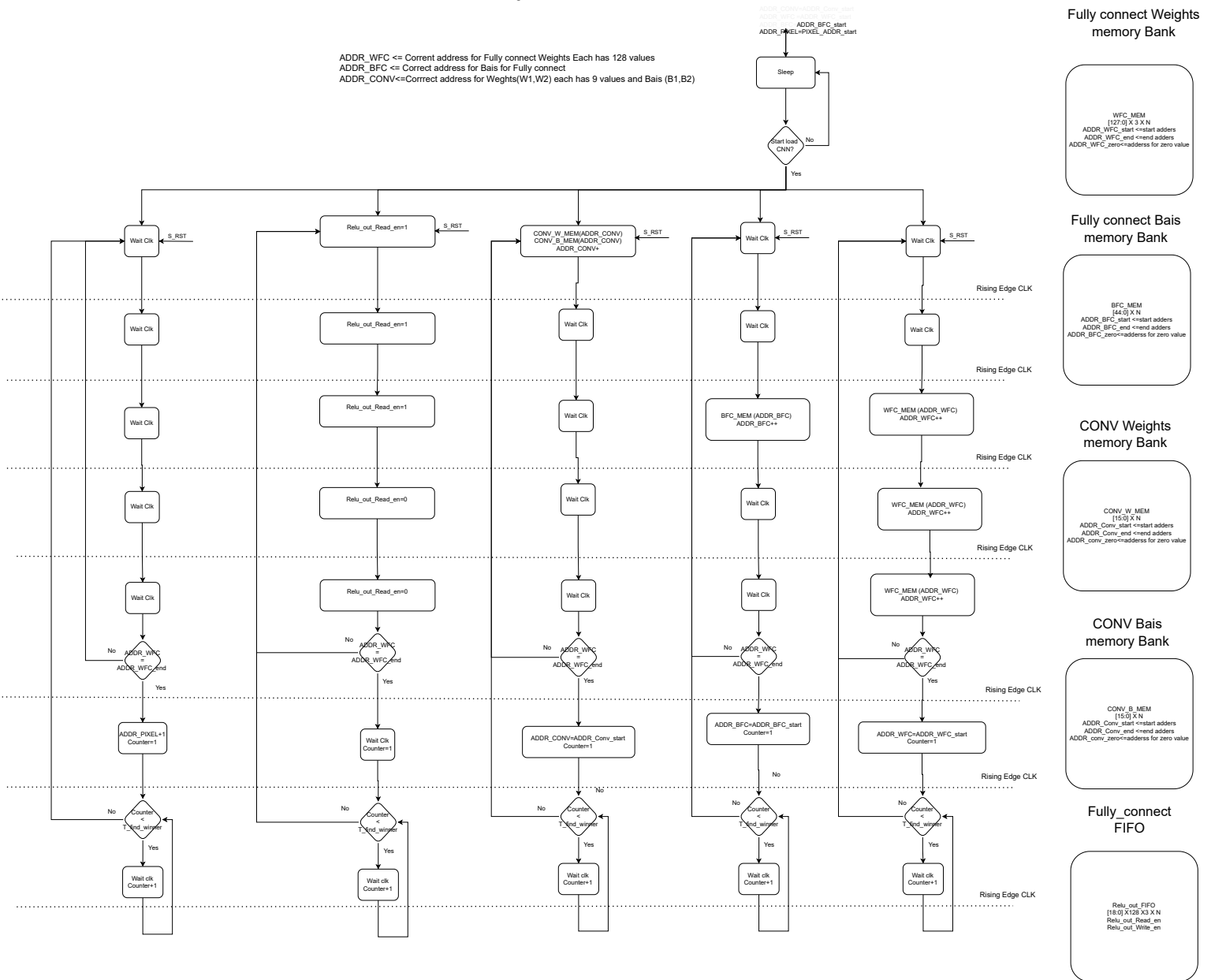


Memory Controller SM

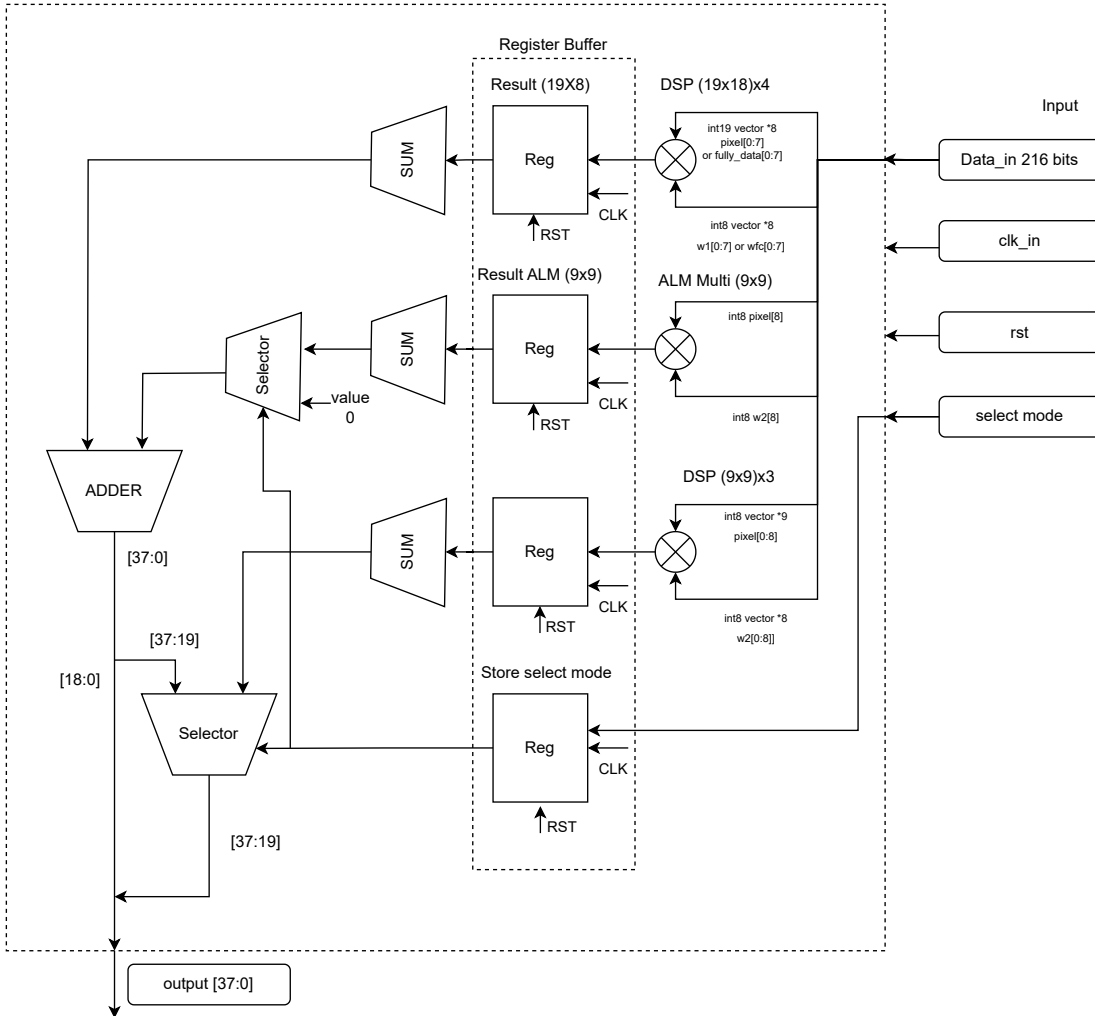


Memory controller work flow

ADDR_WFC <= Correct address for Fully connect Weights Each has 128 values
ADDR_BFC <= Correct address for Bais for Fully connect
ADDR_CONV<=Correct address for Weights(W1,W2) each has 9 values and Bais (B1,B2)



MultiMultiplierEngine (created and tested)



Block operation:

When Mode 0 selected:

We perform Sum of
 $\text{eight}(\text{data}(\text{int19}) * \text{weight}(\text{int8}))$
 Input split :
 $\text{pixel}_{\text{int8}}[0:8]$ 72 bits
 $\text{weight1}_{\text{int8}}[0:8]$ 72 bits
 $\text{weight2}_{\text{int8}}[0:8]$ 72 bits

When Mode 1 selected:

We perform Sum of
 $\text{Nine}(\text{pixel}(\text{int8}) * \text{weight1}(\text{int8}))$
 and the sum of
 $\text{Nine}(\text{pixel}(\text{int8}) * \text{weight2}(\text{int8}))$
 Input split :
 $\text{fully_data}_{\text{int19}}[0:7]$ 152 bits
 $\text{fully_weight}_{\text{int8}}[0:7]$ 64 bits

7 DSP per block:

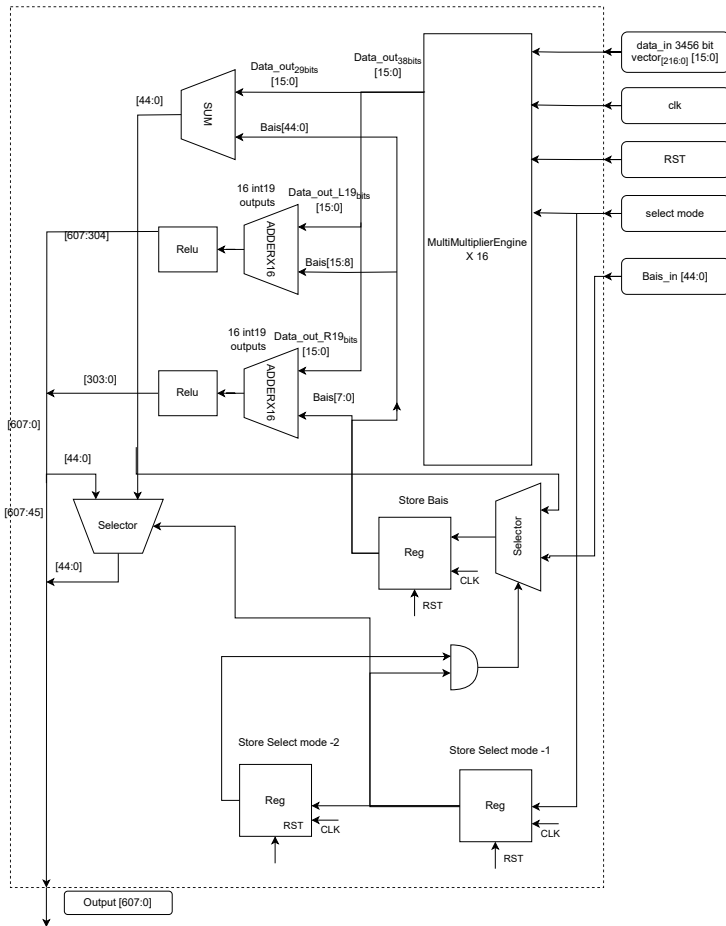
3 configure to run tree 9x9
 multilper each

4 configure to run two 19x8
 multilper each

Pipline configuration

The Block can run with speed up to
 73Mhz without Pipline up to 51Mhz

Parallel_Compute_Engine_16



Block operation:

First we run all the data_in vector in 16 MultiMultiplierEngine units.

When Mode 0 selected:

We add the biases to the outputs of the 16 MultiMultiplierEngine units for total of 32 different int19 results also doing this step we add selector that will set zero to the negative (Relu layer).

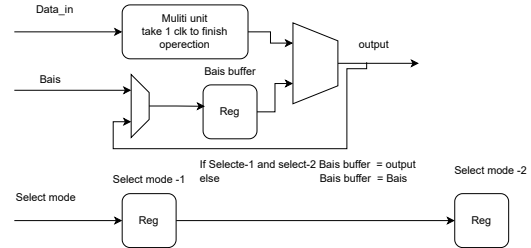
```
Tensor output =
```

Input split :
Bias for the 1

```
Bais for the left results <=Bais_in [7:0]
```

When Mode 1 selected:

We sum the outputs of the 16 MultiMultiplierEngine units with the input bias as it is.



Pipeline								
Clk	<u>Data in</u>	<u>Bais In</u>	<u>Select mode</u>	<u>Mult unit Buffer</u>	<u>Select mode -1</u>	<u>Select mode -2</u>	<u>Bais Buffer</u>	<u>Output</u>
1	Conv(1)	B1,B2	0					
2	0	0	1	Multi Result(Conv(1))	0		B1,B2	
3	date for(383,256)	Bais	1	0	1	0	0	Res(191:175,15:0)
4	date for(255,128)	Dont care	1	Multi Result(date for(383,256))	1	1	Bais	0
5	date for(127,0)	Dont care	1	Multi Result(date for(255,128))	1	1	Sum(383:256)+bais	Sum(383:256)+bais
6	Conv(2)	B1,B2	0	Multi Result(date for(127,0))	1	1	Sum(383:256)+Sum(255:128)+bais	Sum(383:256)+Sum(255:128)+bais
7	0	1	1	Multi Result(Conv(1))	0	1	B1,B2	Fully correct result