

Design of a parameterized CNN accelerator

Cheryl (Yingqiu) Cao

2022-08-06

Goal

To design and to implement a systolic array of MACs, double buffers for the ifmap, weights and ofmap, address generation logic for each buffer, and the top level state machine that orchestrates the flow of data to produce desired outputs.

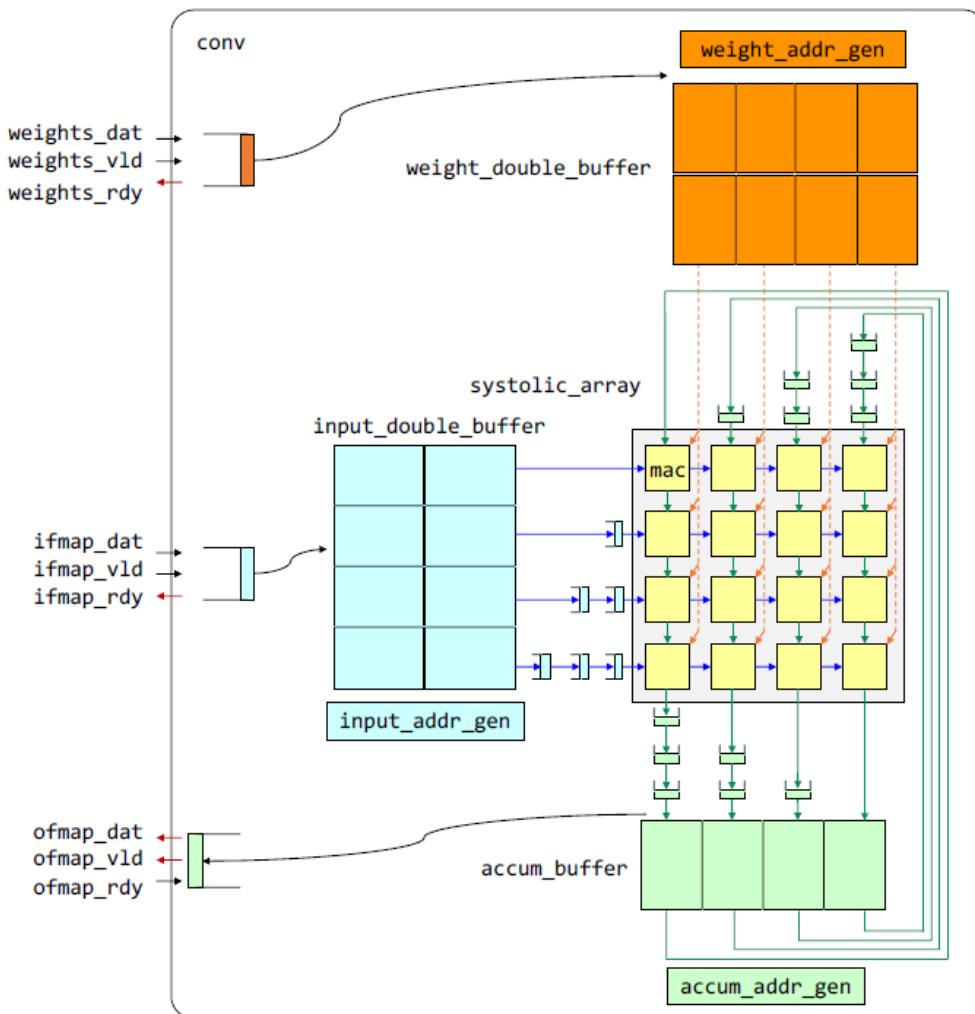


Figure 1: CNN accelerator block diagram.

Analysis

Design

Sync of MAC array (main FSM)

Mac array

ifmap input control

input/weight chaining module

double buffer module

ofmap output control

accum double buffer

ofmap_PISO

Tools

Project Files

Verilog CNN implementation files

Testbench files

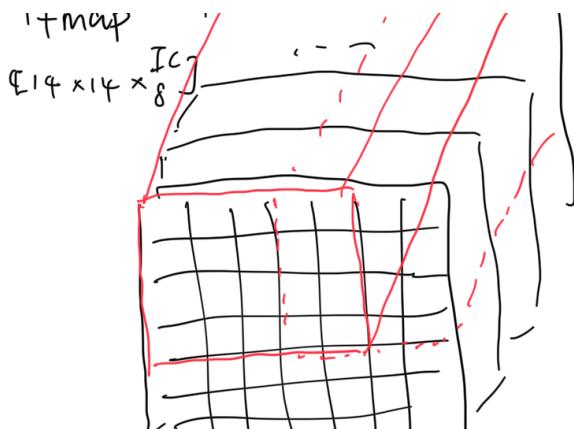
Typedef files

Analysis

input double buffer

ifmap buffer saves IC1 = 2, IY0 = 5, IX0 = 5, IC0 = 4 in each bank

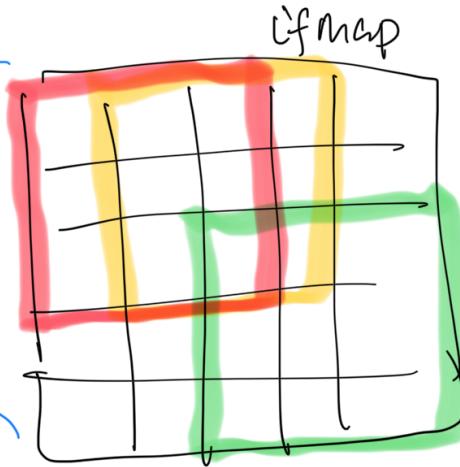
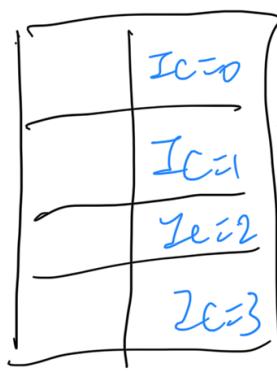
There are a total of OY1 * OX1 * OC1 = 4 * 4 * 4 = 64 banks



double layer buffer

- loads ifmap [5x5x8] in each bank
- next bank overlaps with this one by 2 columns

input channel buffer

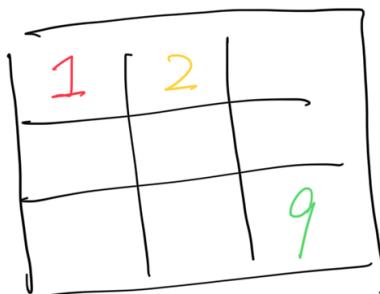


○ 框的大小是 ofmap
3x3 - 三分之二

feature map

$w(:, :, IC=0, OC=x)$

○ 框的个数
feature map
size - 三分之二



为了得到 3×3 的 ofmap,
只 red 框里的 ifmap data
需要 convolve $w(0, 0, :, :)$

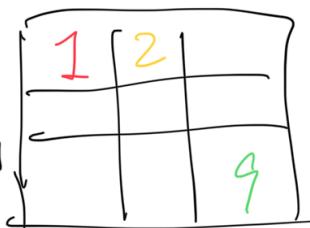
input buffer 扩出到 MAC array

的顺序是

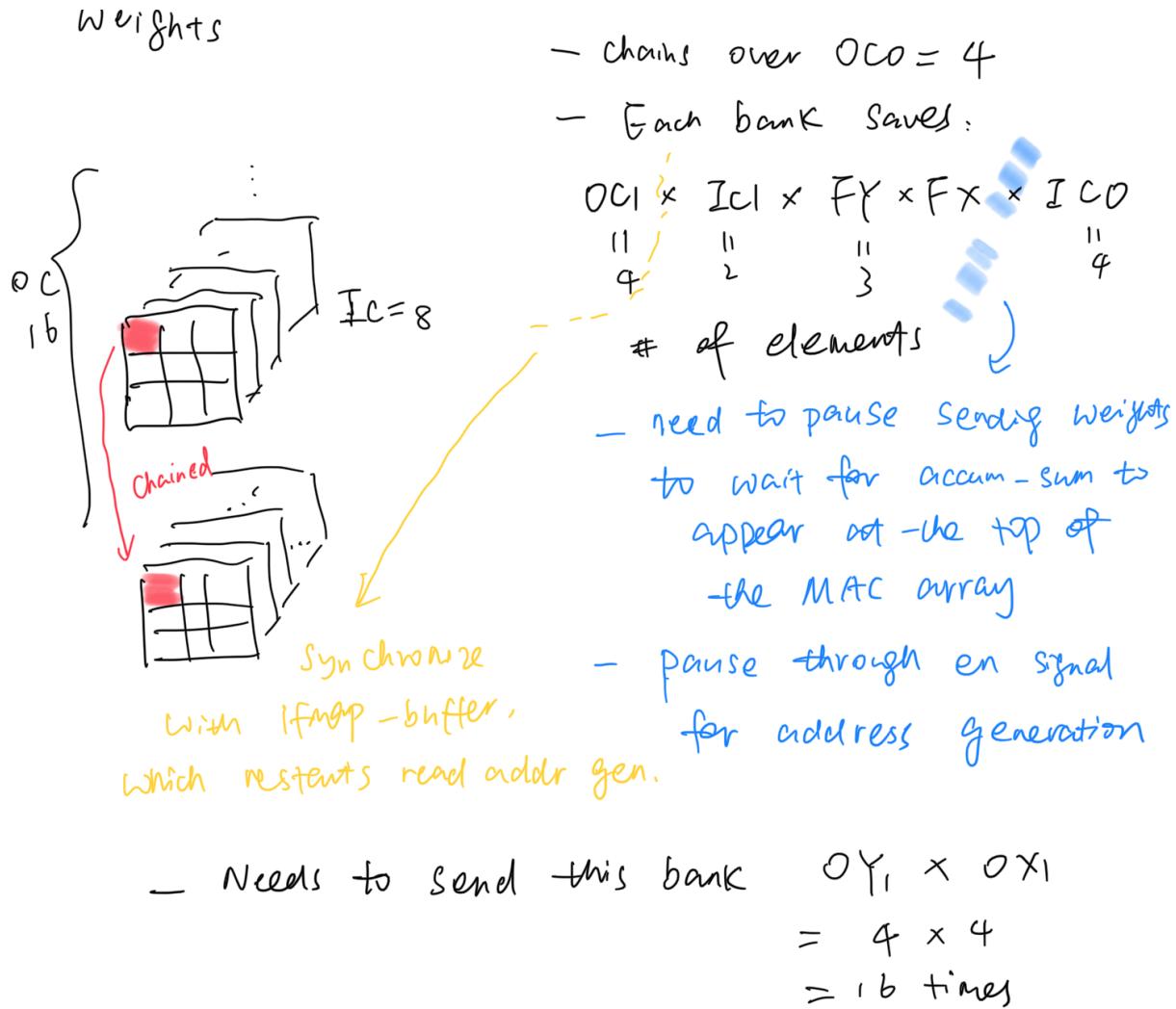
红框 $I(0,0,0) [1,0,0] (2,0,0)$
 $(0,1,0) [1,1,0] (2,1,0)$
 \downarrow $[2,2,0]$

黄 \rightarrow 绿 $I(2,2,0) [3,2,0], (4,2,0)$
 $(4,4,0)$

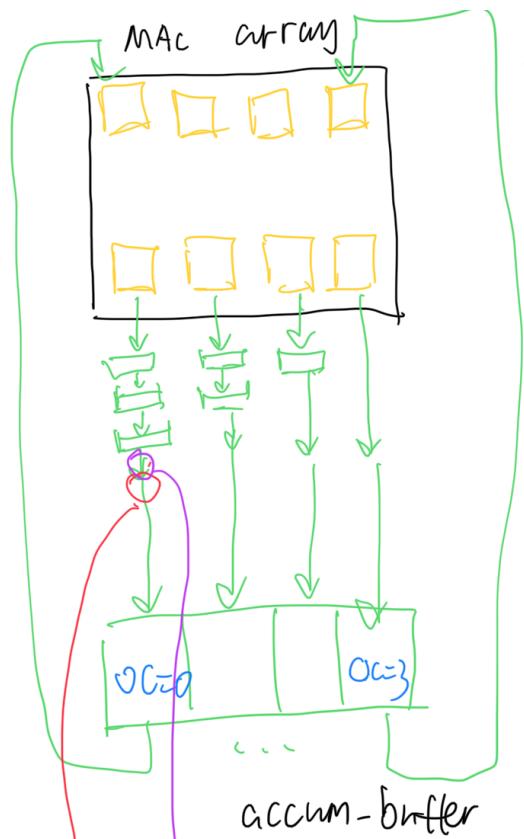
ofmap



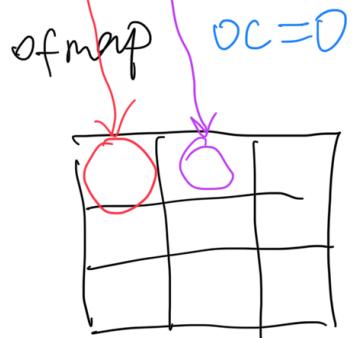
Weight double buffer



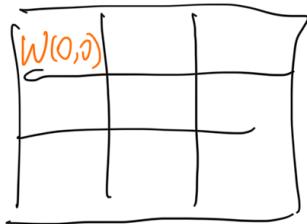
Accum_buffer



- △ first output is a partition
Sum of $O(0,0, OC=0)$
for $IC=0:3$
- △ 2nd output is a partition
Sum of $O(1,0, OC=0)$
for $IC=0:3$
- △ After 9 (ofmap size) cycles,
all elements in ofmap Convolved
with $W(0,0, IC=0:3, OC=0)$
- △ "feed"
accum-buffer results back
to MAC array.



$OC=0$
feature map
 $IC=0:3$



The 1st output is $O(0,0, oc = 0)$ for $w(0,0)$, $IC = 0:3$

The 2nd output is $O(1,0, oc = 0)$ for $W(0,0)$.

...

The 9th output (9 = ofmap size)
is $O(2,2, oc = 0)$ for $W(0,0)$

The 10th output is $O(0,0, \text{oc} = 0)$ for $W(1,0)$ -- the 1st output needs to be added to this from the `accum_buffer`.

Once we iterate through one bank of weight buffer, we are done for $O(0:2, 0:2, \text{oc} = 0:3)$ for $\text{IC} = 0:3$.

- next step is to feed Ofmap back through the `accum_buffer`.

- The weight bank will cover $W(:, :, \text{oc} = 0:3)$ for $\text{IC} = 4:7$.

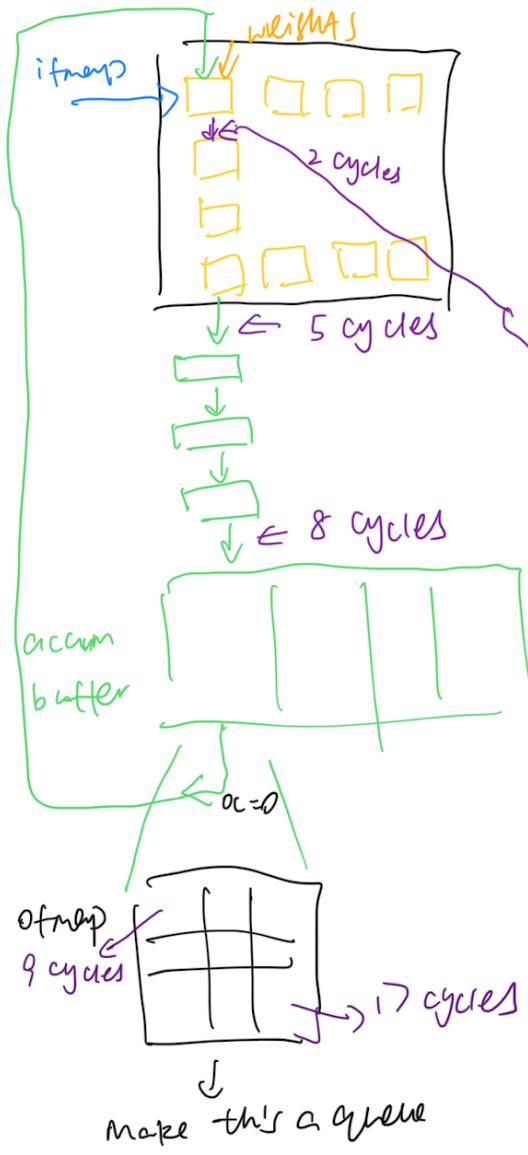
After that we need to reset `accum_buffer` because OC will move to 4:7.

After we iterate through all weight banks in the weight buffer, one blocking of ofmap is done -- $O(0:3, 0:3, \text{oc} = \text{all})$.

DRAM then reads the next block from the original ifmap.

Synchronization of the MAC array and the double buffers

- See the figure below for general analysis



- $t=-1$ - read weight-addr-gen
- $t=0$ - weights appears at the top
enable ifmap-addr-gen
 $w(0,0, t=0)$
- $t=1$ - ifmap(0,0, OC=0, CC=0)
appears $F_x=0, F_y=0$
- $t=2$ - $w(0,0, CC=2)$ appears
partial sum for
 $O(0,0, OC=0, CC=0)$ appears here
- $t=8$ - partial sum appears at the
// input of accum-buffer
- enable addr-gen for
accum-buffer
- - - - -
 $t=9$ - $w(1,0, CC=0)$ appears
- com degre - ifmap(2,2, OC=0) at input
- enable ifmap
- $t=10$ - partial sum appears at
the top of MAC array
- ifmap(0,0, OC=0, CC=0)
 $F_x=1, F_y=0$
appears
- element
- $\leftarrow \boxed{1} \boxed{(1,0)} \boxed{(0,1)} \dots \boxed{1} \rightarrow \text{enqueue}$

$O(0,0)$
- accum-but 1st element ready to degre w/ 9 cycles

- accum-cycle = $I_{CO} + O_{CD} + 1 = 9$

- ifmap iteration over 1 weight takes 9 cycles

ifmap-cycle = $O_{Y0} \times O_{X0} = 9$

- if ifmap-cycle is larger (\geq) $T = \max(\text{accum-cycle}, \text{ifmap-cycle})$

- pauses weights longer for $T - I_{CO}$

- if accum-but-cycle is larger:
 $= 9 - 4 = 5$ cycles

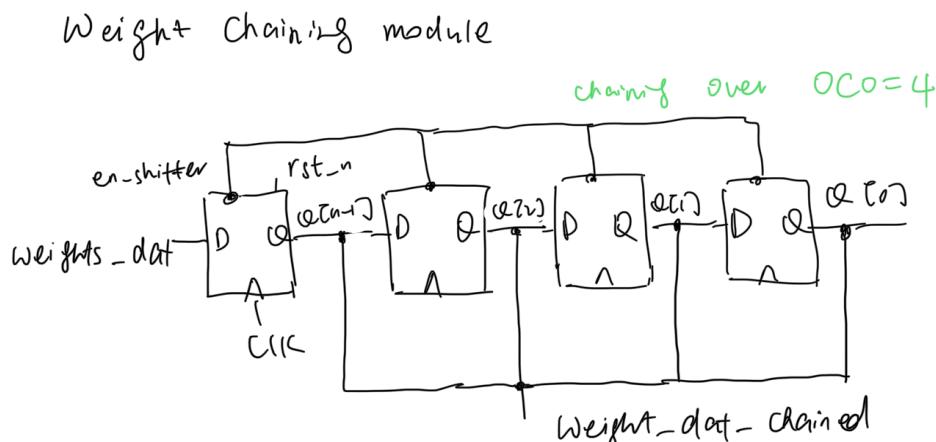
- pauses ifmap-read-addr-gen for $T - \text{ifmap-cycle}$

- Synchronization case 2: ifmap_cycle >= accum_cycle
 - It takes longer to iterate through the ifmap elements for weight ($F_x = 0, F_y = 0$), than for the 1st partial sum ofmap(0,0, IC = 0:3) to get out of the MAC array
 - need to wait to feed back the accum_buf to MAC array

Design

input/weight chaining module

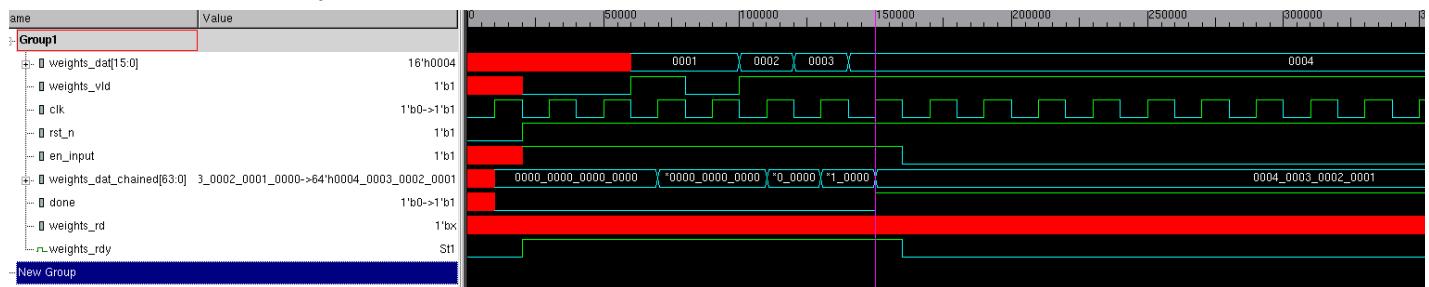
- takes ifmap_data from the testbench, chains them along IC0 = 4, outputs them to the input_double_buffer
- takes weights_data from the testbench, chains them along OC0 = 4, outputs them to the weights_double buffer
- The following figure shows the design for the weight chaining module.
 - for the input changing, replace “weights” with “input”, replace “OC” with “IC”.



```

local vars:
interface
  readif = weights-vid_bb-ready;
  en-shifter = en-input && (input[75:0] == weights-data);
  weights-chaining # (parameter OC0=4)
    ( input[75:0] weights-data,
      weights-data);
  end
end
  
```

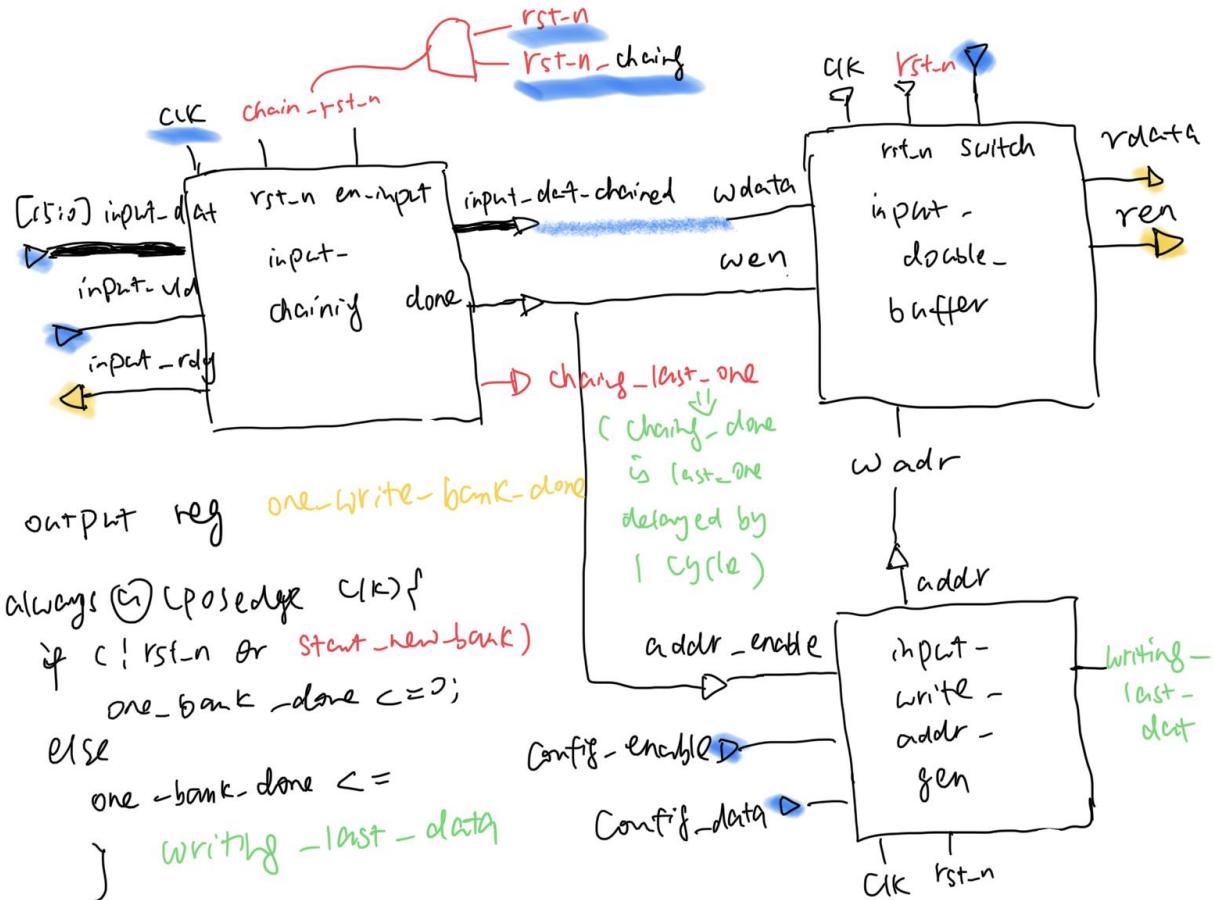
- The simulation waveform is as below
 - “done” goes high when chained_data appears at the output, and stays high until cleared by the state-machine/controller.



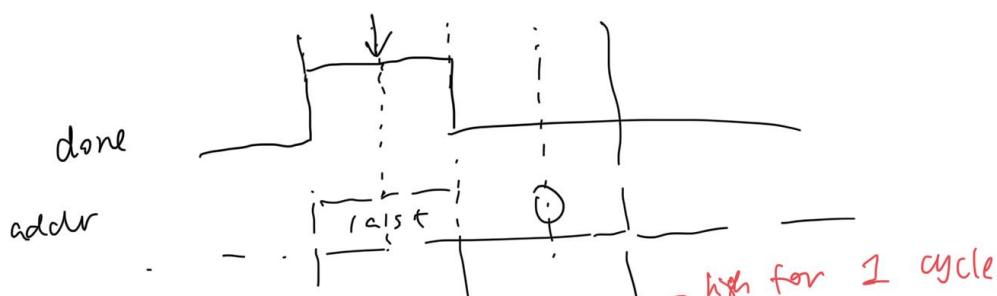
ifmap input control

- Takes ifmap data from the test bench
- sends ifmap data from the double buffer to the MAC array
- write_bank_count tracks the # of banks we have finished writing to.

Input Control



- we want "done" to stay high for only 1 cycle
- therefore, need to: reset "chaining" module immediately
 - to clear "done" flag
 - to clear SIFO
 - to block additional input



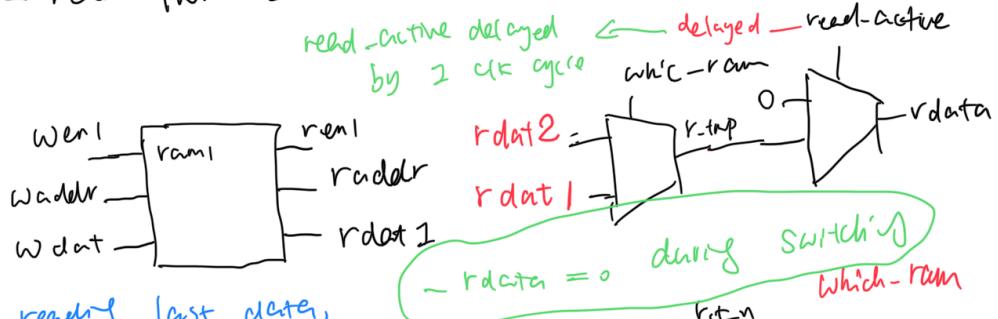


logic read/write-bank-ready-to-Switch;
logic ready-to-Switch;

double buffer module

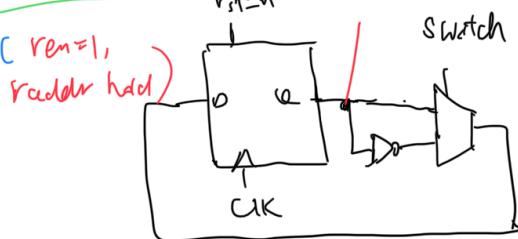
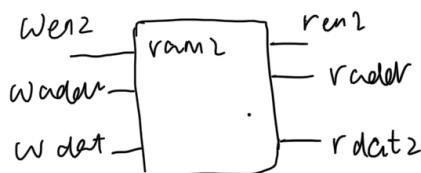
double buffer

- read from 1 ram, write to the other



- After reading last data,

pause 1 cycle before switch



Local Signals:

```

logic which_ram;
logic active;
logic write_active;
logic read_active;
logic [DataWidth-1:0] rdat1,
logic [DataWidth-1:0] rdat2;
active = rst-n && (!switch),
write_active = active && wen;
read_active = active && ren;
wen1 = write_active && (which_ram)
wen2 = write_active && (!which_ram)
ren1 = ...

```

Diagram showing the local signals for the double buffer module. The signals include which_ram, active, write_active, read_active, rdat1, rdat2, and active. The active signal is defined as rst-n && (!switch). The write_active signal is defined as active && wen. The read_active signal is defined as active && ren. The wen1 signal is defined as write_active && (which_ram). The wen2 signal is defined as write_active && (!which_ram). The ren1 signal is defined as A note indicates that must rst the double buffer before use and initialize switch-bank to 0 during reset.

*Notes:

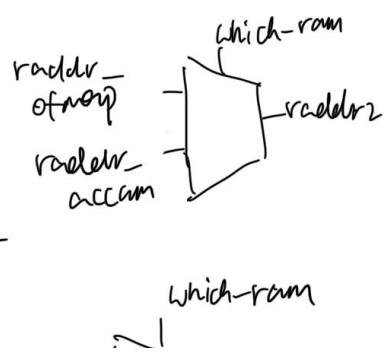
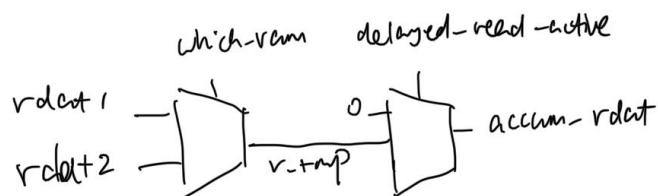
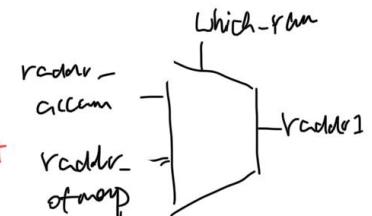
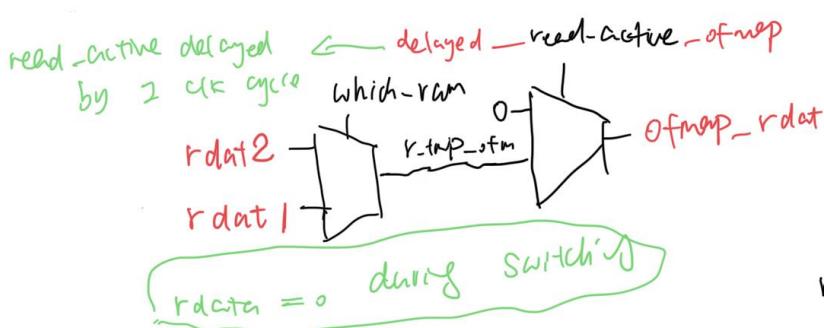
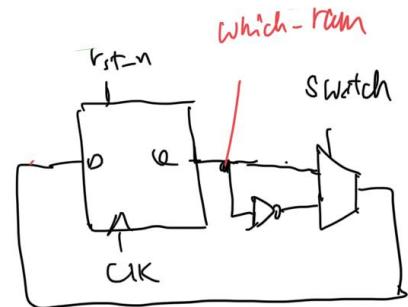
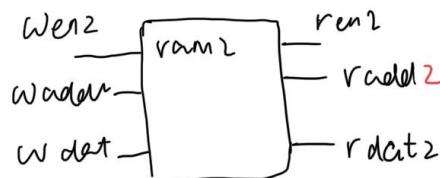
- must reset the double buffer before use
- during reset, initiate switch_banks to 0
- after reading the last data, keep ren =1 w/ the same raddr for another cycle before pulling switch high

accum double buffer

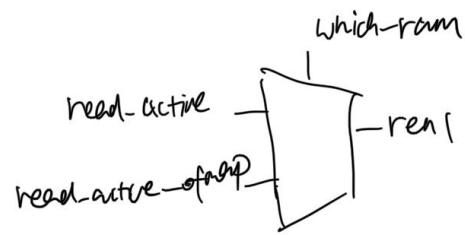
- different from the general double buffer used for ifmap and weight data, accum double buffer has 1 write port but **2 read ports**.
- one read port feeds ofmap partial sum back to the MAC array. The other one sends completed ofmap data out to the testbench.



- After reading last data,
pause 1 cycle before switch



$rdat2 \rightarrow [rtmp]$



local signals:

```

logic which-ram;
logic active;
logic write-active;
logic read-active, read-active-ofmap;
logic wen1, wen2, ren1, ren2;
logic [DataWidth : 0] rdat2;

```

active = $\overline{rst-n} \wedge \overline{(\text{!switch})}$;

$$\text{write-active} = \text{active} \wedge \text{wen};$$

$$\text{read-active} = \text{active} \wedge \text{ren};$$

$$\text{read-active-ofmap} = \text{active} \wedge \text{ren-ofmap};$$

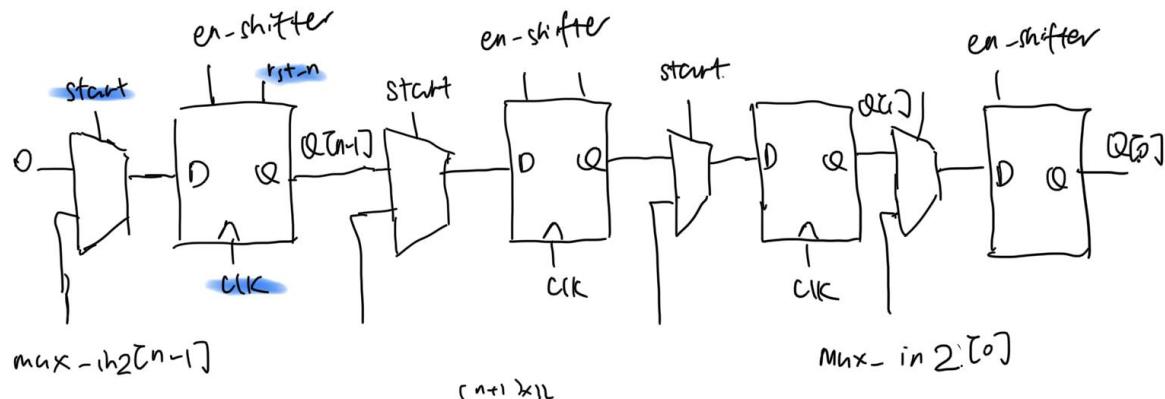
$$\text{wen1} = \text{write-active} \wedge \overline{(\text{!which-ram})}$$

$$\text{wen2} = \text{write-active} \wedge (\text{which-ram})$$

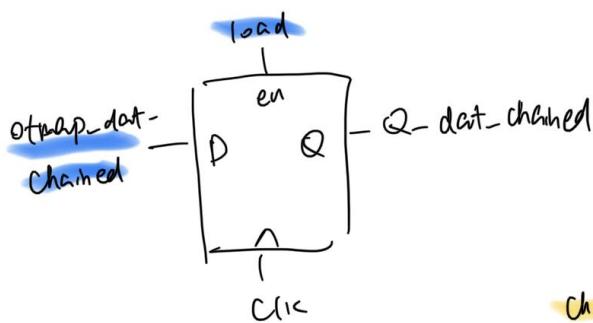
ofmap_PISO

ofmap_PISO

$n = 0 \text{CO}$



$$\text{mux_in}_2[n] = Q_{\text{dat-chained}}[-1 : n \geq 2] \quad \text{ready_to_unchain} = \text{ofmap_rdy};$$



en-shifter = en-PISO &

ready-to-unchain;

- use en-shifter as en signal
for unchain-Counter;

chaining-last-one = en-shifter & (counter ==
 \checkmark $(n-1)$).

local signals

interface

ofmap_PISO # (parameter OCO = 4)

C	input	clk,
	input	rst-n,
-	input	en-PISO,
	input	load,
	input	start,
	output	chaining-last-one,

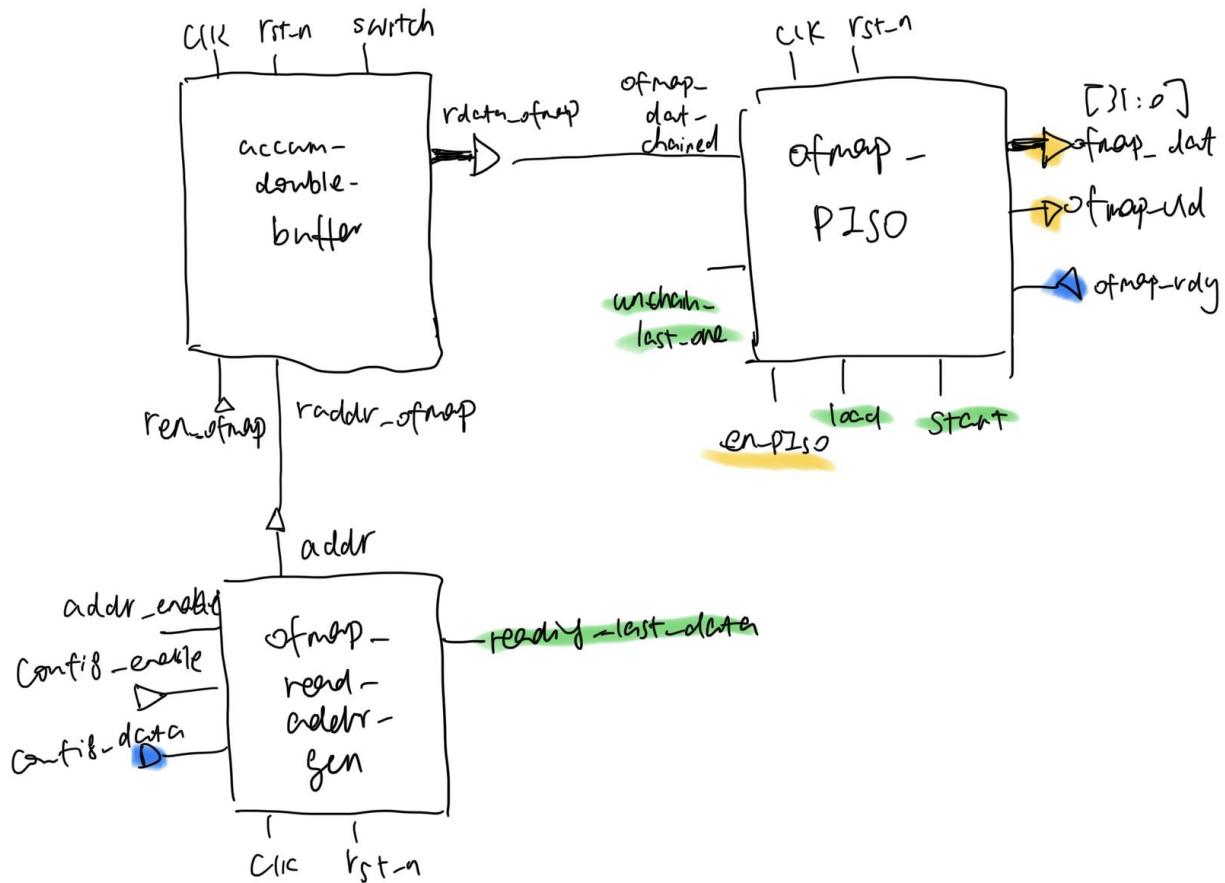
input [32 * OCO - 1 : 0] ofmap_dat-chained,
input ofmap_rdy,
output [32 : 0] ofmap_dat,
out put ofmap_vld);

only 1 for one cycle per
unchained data

ofmap output control

```
+1 ofmap_FSM.v | 3 ifmap_input_FSM.v
4 // Author: Cheryl (Yingqiu) Cao
5 // Date: 2022-03-12
6
7 module ofmap_FSM
8 (
9   input logic clk,
10  input logic rst_n,
11
12  // for the read_addr_gen
13  input logic last_ofmap_data,      // the unchained data in the last one in this ofmap bank
14  output logic config_enable,
15  output logic raddr_gen_en,
16
17  // for accum_double_buffer
18  output logic ren,
19  output logic siwtch,
20
21  // for ofmap unchaining PISO
22  input logic unchaining_last_one,
23  input logic ready_to_unchain,
24  output logic en_PISO,
25  output logic load,
26  output logic start,
27
28  // for the ofmap read bank counter
29  output logic one_read_bank_done,  // en signal for the counter
30
31  // for the main FSM
32  input logic ready_to_switch,
33  input logic start_new_read_bank,
34  output logic read_bank_ready_to_switch
35
36 );
37
```

ofmap-output-Controller

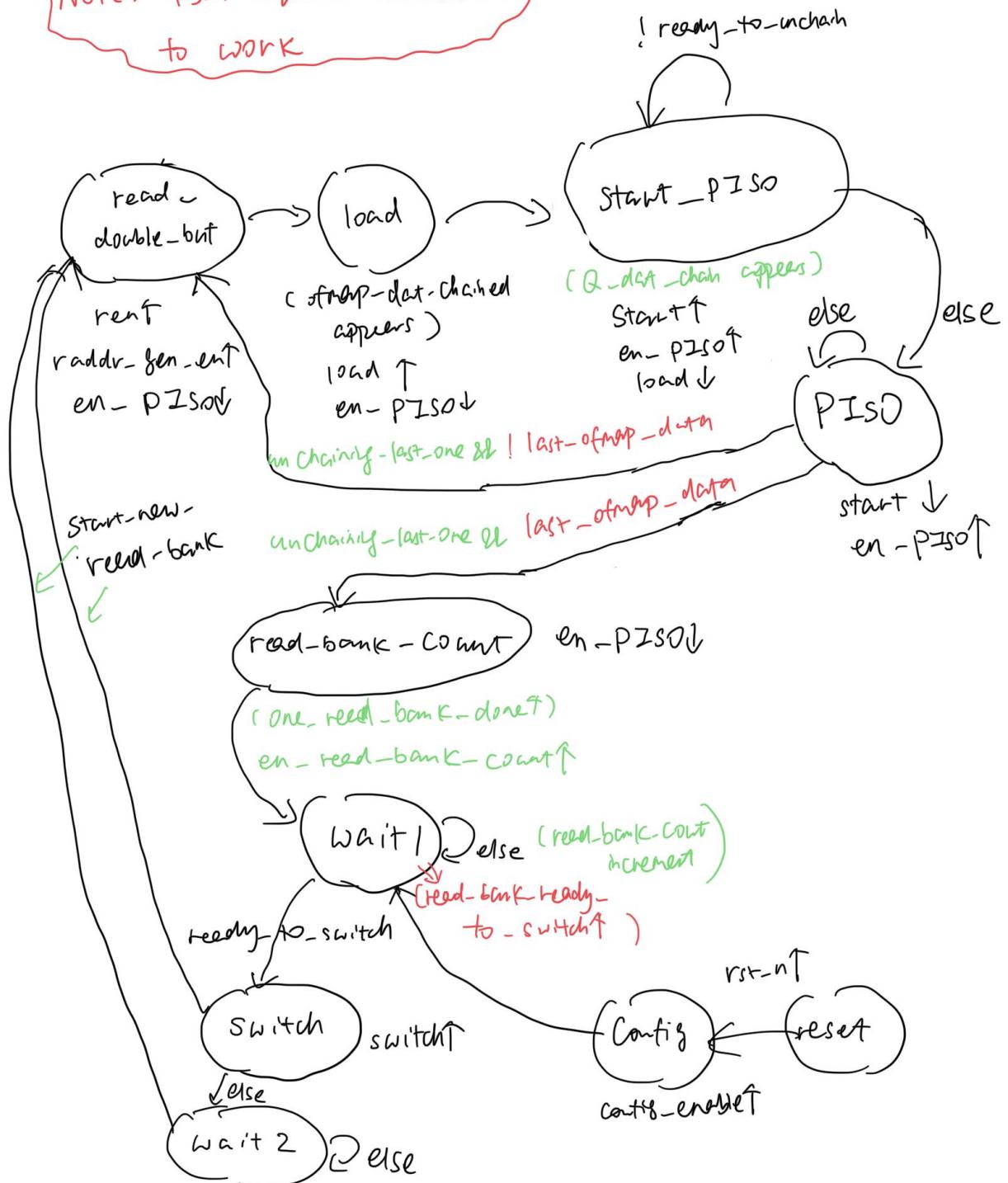


```

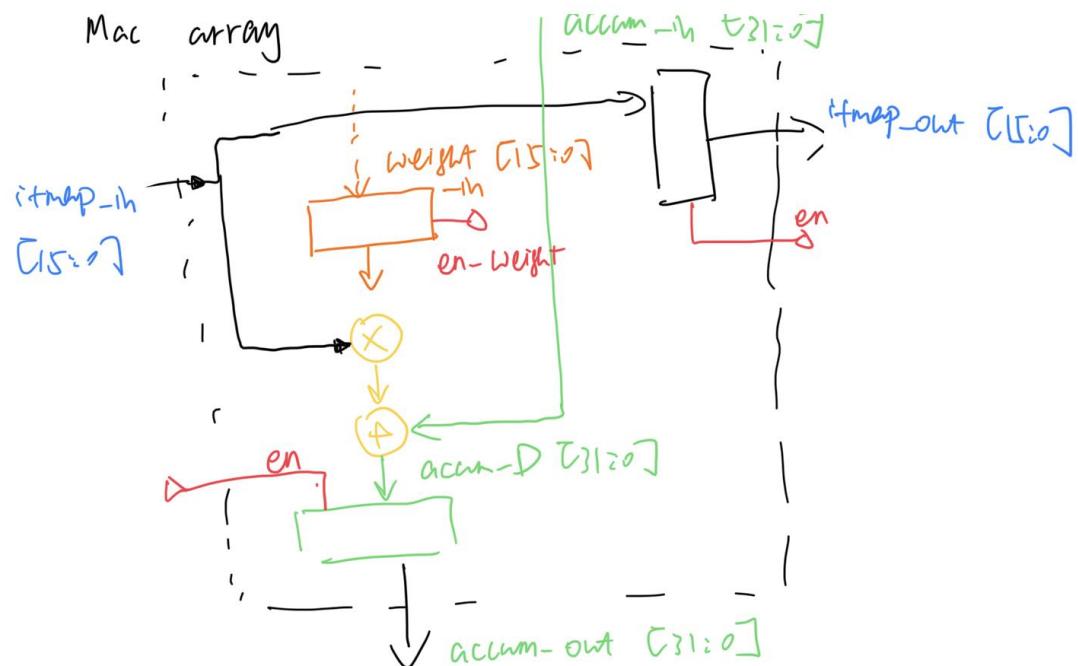
always @(posedge clk) begin
    if (rst || switch)
        last_of_msp_data <= 1;
    else if (ready & !last_data)
        last_of_msp_data <= 1;
end

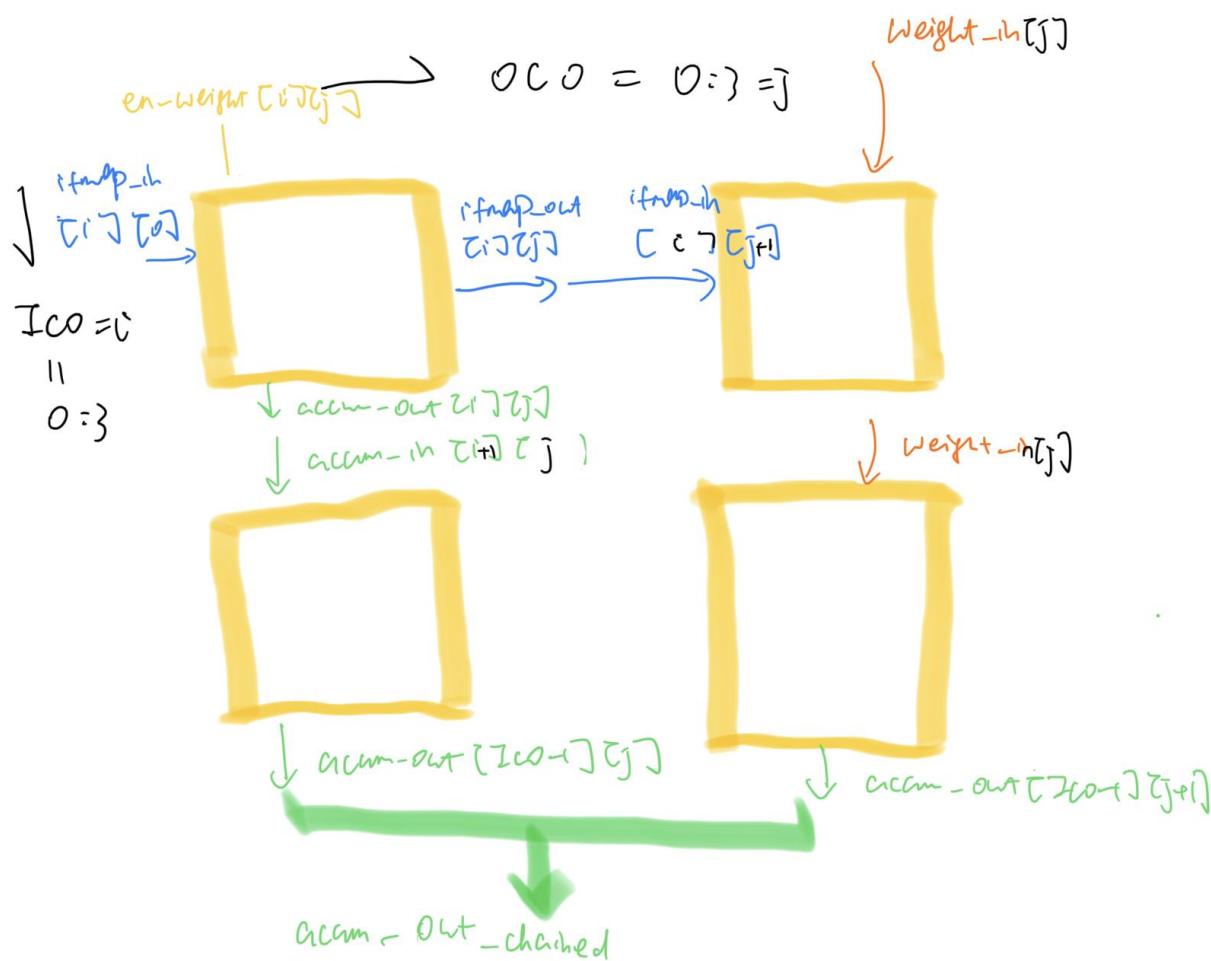
```

Note: FSM requires $n=OCO > 1$
to work



Mac array





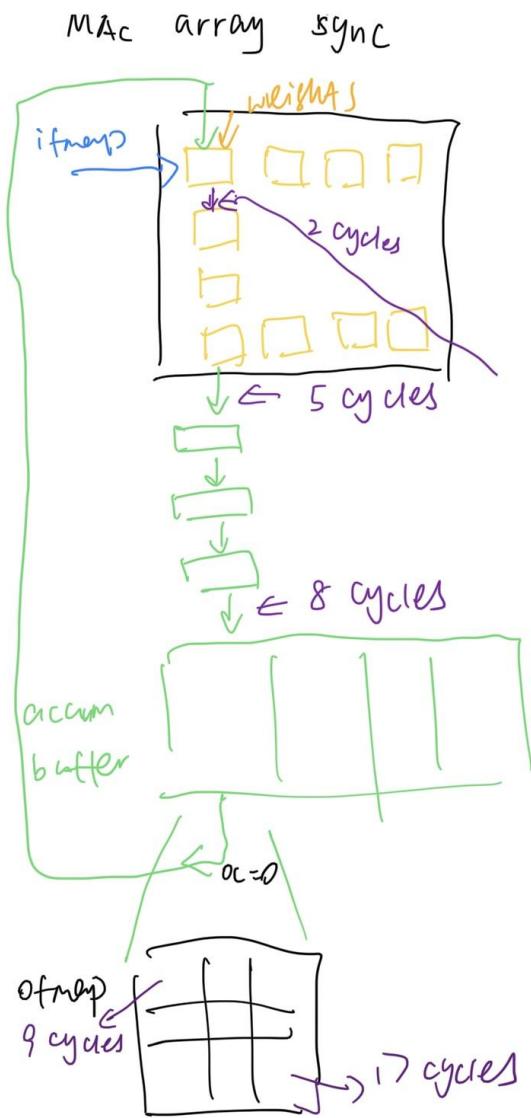
$$ifmap_in[i][j] = ifmap_out_chained[16(i+1)-1 : 16i]$$

$$weight_in[j] = weight_out_chained[16(j+1)-1 : 16j]$$

$$accum_in[0][j] = accum_out_chained[32(j+1)-1 : 32j]$$

$$accum_out[ICO-1][j] = accum_out_chained[32(j+1)-1 : 32j]$$

Sync of MAC array



$t=0$ weight-double-buffer: $\text{ren} \uparrow$
 ifmap-double-buffer: $\text{ren} \downarrow$

$t=1$ (weight appears at top of mac)
 $w[0, 0, \overset{\uparrow}{\text{icoz}}]$

weight-fifo-eng, ↑
ifmap-fifo-eng ↓
mac-more.: em-weight 00 ↑

$t=2$ [ifnewp / accum_m data
appears (in) top/left of Mac)

weight-fifo-eng ↑
 max-more: en-Weight OOK
 if map-fifo-eng ↑
 accum-in-fifo-eng
 C = (~1st-FKFY-cycle)

weight-double-buffer: rent
ifmap-double-buffer: rent↑
accum-double-buffer:
 $\leftarrow c \sim \text{1st-FIFO-cycle} \right)$

$t = \text{ICO}$ (pause weight ready)

weight_double-buffer: rend

$$t = ICO + 1$$

weight-fifo-enq ↓

mac-more: en-weight > 0 ↓

accum-out-fifo-enq ↓

$$t = ICO + 2$$

(accum-out appears @ top

of fifo)

accum-out-fifo-enq ↑

$$t = ICO + OCO + 1 = 9$$

(accum-out appears @ top of
double-buffer)

accum-double-buffer = wen ↑

} assign en-ic1-fy-fx-counter = (OY0-OX0 == Max);

} assign ic1-fy-fx-iter-done = (IC1-fy-fx == Max) &

} assign OCL-iter-done = (OCL == (OCL-1)) & en-OCL-counter;

assign en- ∂Y_1-OX_1 -counter = OCL-iter-done

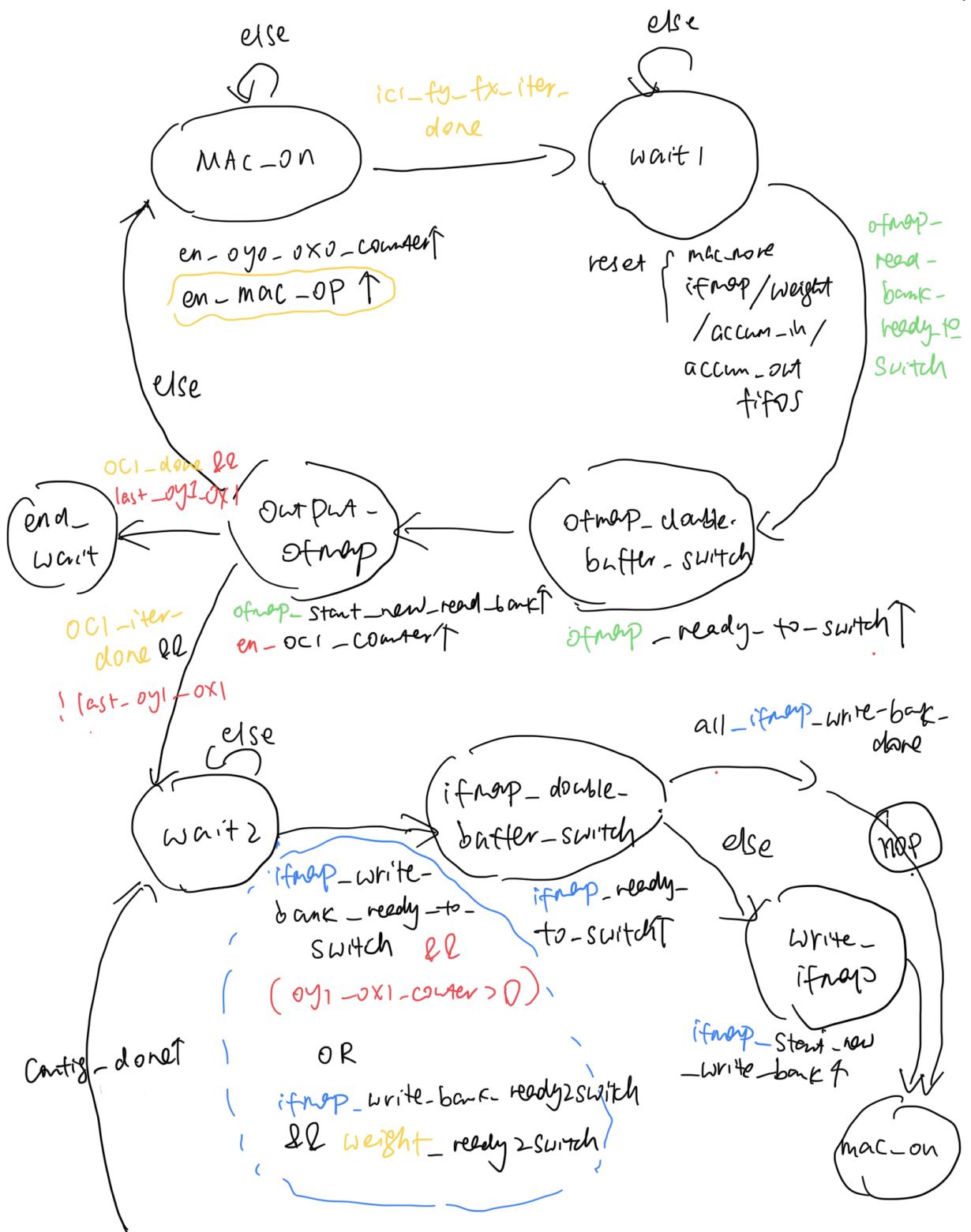
for loop

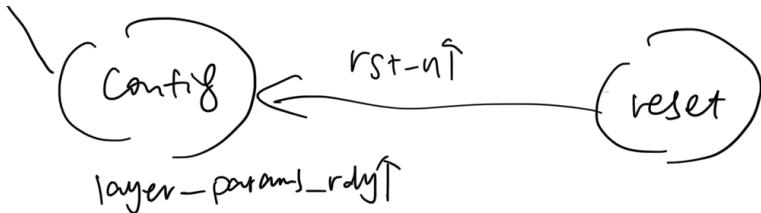
assign all-ifnop-write-bank-done = (write-bank-count

assign weight-ready-to-switch = == ∂Y_1-OX_1);

↓ ifnop-ready-to-Switch

only 1 weight bank & $(\partial Y_1-OX_1\text{-counter} == 0)$;





- Assuming $OY0 * OX0 \geq OC0 + IC0 + 1$
- oy0_ox0, a counter that increments each clk cycle
 - range is [0, OY0_OX0]
 - Use $OY0_OX0 + 1 = 3*3 + 1 = 10$ as period
- ic1_fy_fx, a counter that increments each OY0_OX0 iteration completes
 - range is [0, IC1_FY_FX], saves the # of OY0_OX0 cycle that was completed

For **mac_more** module:

Control signal	en_weight00	weight_fifo_enq	ifmap_fifo_enq	accum_in_fifo_enq		accum_out_fifo_enq	
On cycle : oy0_ox0	[1,1] oy0_ox0_first_cycle	[1,IC0]	[2, OY0_OX0]	[0,0]	[2, OY0_OX0]	[0,0]	[IC0 + 2, OY0_OX0] [0, IC0]
Condition: ic1_fy_fx	[0, IC1_FY_FX-1] ic1_fy_fx_not_last_cycle	[0, IC1_FY_FX-1] ic1_fy_fx_not_zero_cycle	>0 ic1_fy_fx_not_zero_cycle	[1, IC1_FY_FX-1] ic1_fy_fx_not_zero_cycle	> 1	[0, IC1_FY_FX-1] ic1_fy_fx_not_last_cycle	>0 ic1_fy_fx_not_zero_cycle

For **weight_double_buffer** module:

Control signal	ren	weight_read_addr_gen_en
On cycle : oy0_ox0	[0, IC0-1]	
Condition: ic1_fy_fx	[0, IC1_FY_FX-1] ic1_fy_fx_not_last_cycle	

For **ifmap_double_buffer** module:

Control signal	ren	ifmap_read_addr_gen_en

On cycle : oy0_ox0	[1, OY0_OX0] oy0_ox0_not_zero_cycle
Condition: ic1_fy_fx	[0, IC1_FY_FX-1] ic1_fy_fx_not_last_cycle

For `accum_double_buffer` module:

Control signal	ren	wen	
On cycle : oy0_ox0	[1, OY0_OX0] oy0_ox0_not_zero_cycle	[IC0+OC0+1, OY0_OX0]	[0, IC0+OC0-1]
Condition: ic1_fy_fx	[1, IC1_FY_FX-1] ic1_fy_fx_not_zero_or_la st_cycle	[0, IC1_FY_FX-1] ic1_fy_fx_not_last_cycle	[1, IC1_FY_FX] ic1_fy_fx_not_zero_cycle

- assign `accum_read_addr_gen_en` = ren;
- assign `accum_write_addr_gen_en` = wen;

testbench ifmap/weight/ofmap data organization (*needs correction)

Test-bench data feeding order

- for ifmap (ix_1, iy_1, ic)

```

for  $oy_1 = 0:3$            one ifmap bank
  for  $ox_1 = 0:3$  -
    for  $ic_1 = 0:1$  -
      for  $iy_0 = 0:5$  -
        for  $ix_0 = 0:5$  chaining
          for  $ic_0 = 0:3$  -

```

- for ofmap (fx, fy, ic, oc)

```

for  $oc_1 = 0:3$  output 1 batch of ofmap data to tb
  for  $ic_1 = 0:1$ 
ofmap read/write
  for  $fy = 0:2$ 
  for  $fx = 0:2$  feeds accum-but
    for  $ic_0 = 0:1$  back to mac array
    for  $oc_0 = 0:3$  chaining
 $oc = oc_0 + oc_1 * (OC_0=4)$   $oc_0=0$   $oc_0=3$ 

```

$$oc = oco + oci * (OC0=4)$$

- for accum_buffer (ox, oy, oc)

for $oy_0 = 0:2$

for $ox_0 = 0:2$

- to re-organize ofmap at tb : (ox, oy, oc)

for $oy_1 = 0:3$

for $ox_1 = 0:3$

for $oci_1 = 0:3$

for $oy_0 = 0:2$

for $ox_0 = 0:2$

for $oco = 0:3$

$$ox = ox_0 + ox_1 \times (OC0=3)$$

Tools

Synopsys VCS for compilation and simulation

Cadence Verdi for waveform visualization and debugging

Synopsys DVE (deprecated by Synopsys during the project)

Commands

```
module unload vcs
```

```
module load vcs/O-2018.09-SP
```

Project Files

Verilog CNN implementation files

Top level

- conv_tiled.v
 - My CNN implementation with MAC array
- conv.sv
 - pseudo CNN implementation
 - directly reads “gold_ofmap” data from file and assigns it to ofmap array

ifmap related

- input_write_addr_gen.v
 - serial port sends data to double layer buffer
 - the addr increases from 0 to $IC1*IX0*IY0$, because **the data for IC0 = 4 channels are chained together**
 - Each bank saves $IC1*IX0*IY0*IC0$ elements of ifmap data
 - **input sequence order** implemented by the test bench:

```
for oy1 in range(0,4):
    for ox1 in range(0,4):
        for ic1 in range(0,2):
            for iy0 in range(0,5):
                for ix0 in range(0,5):
                    for ic0 in range(0,4):

                        ix = ix0 + ox1 * FX;
                        iy = iy0 + oy1 * FY;
                        ic = ic1 * ic0
                        ifmap( ix, iy, ic)
```



- input_read_addr_gen.v
 - it generates the read address for the input double buffer → **MAC array reads data from the ifmap double buffer**
 - it implements the big for loop in the double buffer level and outputs addr in that order for (IC1, FY, FX, OY0, OX0)
 - this is the iteration required for sum accumulation
 - We **still need to iterate through OC2** for blocking
 - **The double layer buffer saves ifmap (IC1, IX0, IY0), and outputs it for OC1 = 4 times for one block**
- input_chaining.v
 - serial-in-parallel-out
 - chains ifmap data over IC0 = 4, before inputting to the ifmap_double_buffer
- ifmap_input_FSM.v
 - controls the timing for input_chaining and writing operations for the ifmap_double_buffer
 - Outputs one_write_bank_done flag to the MAIN FSM
 - takes an input start_new_write_bank from the MAIN FSM to start the next bank
- ifmap_input_controller.v
 - connects ifmap_chaining, ifmap_double_buffer, input_write_addr_gen and other small logic blocks

Weights related

- weights_chaining.v
 - serial-in-parallel-out
 - chains weights data over OC0 = 4, before inputting to the weight_double_buffer
- weight_addr_gen.v
 - addr generation for the read/write operation of weight double buffer
- weight_input_controller.v
 - connects weight_input_chaining, weight_double_buffer, weight_write_addr_gen and other small logic blocks

Ofmap related

- `ofmap_PISO.v`
 - unchains ofmap data over OC0 = 4, before outputting to the serial bus of the test bench.
- `accum_addr_gen.v`
 - generates read/write address for the `accum_double_buffer`
 - goes from 0 to ($OY0^*OX0 - 1$)
- `accum_double_buffer.v`
 - Double buffer for the accumulated sum and ofmap data
 - Different from the general double buffer used for ifmap and weights, this buffer has 1 write port but 2 read ports.
 - One read port and one write port uses the same bank to save the partial accum sum. The other port is used to read ofmap data output for the testbench.

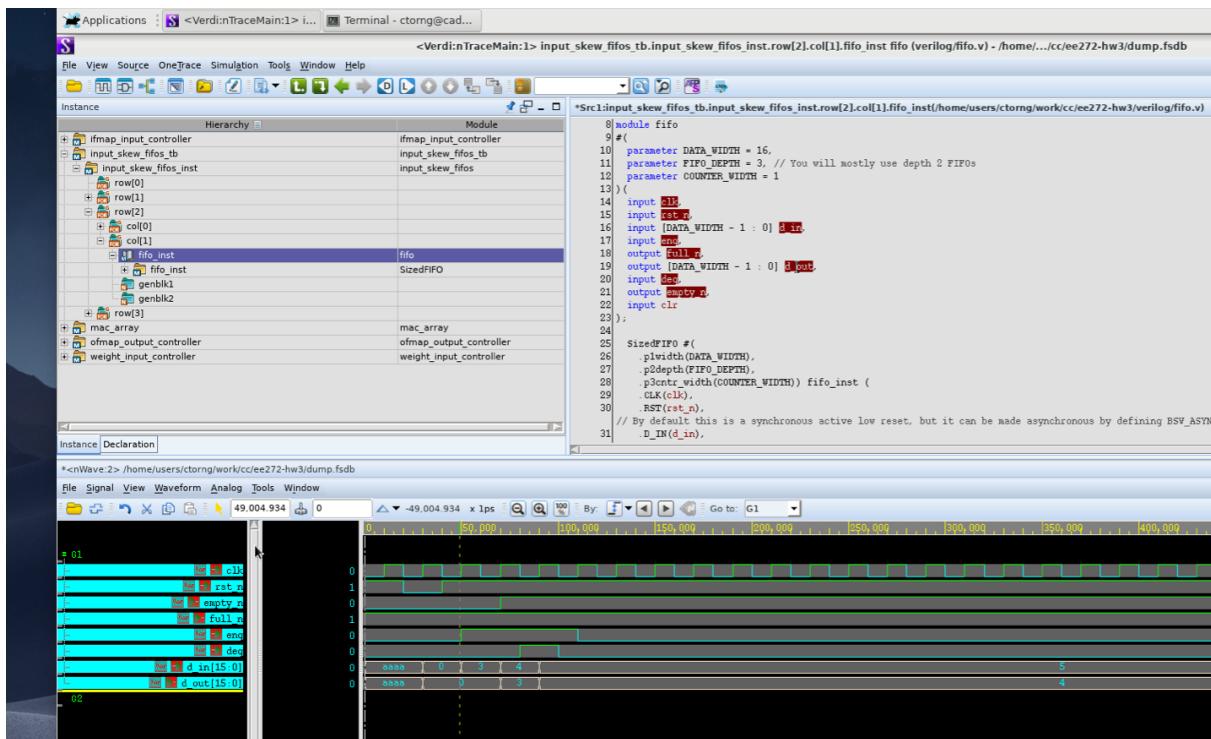
MAC related

- `mac.v`
 - the mac cell
- `mac_array.v`
 - mac array for IC0 = 4 and OC0 = 4
- `en_weight_shifter.v`
 - derives the `en_weight[ic0][oc0]` signal for the `mac_array` from a single `en_weight` signal for the 1st mac cell [`icc = 0`][`oc0 = 0`].

Memory related

- `double_buffer.v`
 - used in two instances to store ifmap and weights locally
 - in each cycle, can perform one read, one write in parallel
- `ram_sync_1r1w.v`
 - can read and/or write one word each clock cycle
 - read data becomes available in the next cycle
- `wrapped_skew_fifo.v`

- **When en = 1, each row dequeues 1 data until the fifo becomes empty**
- add an en signal: enqueue and dequeue stays 0 unless en = 1;
- one enq signal for the entire module: input arrives at the same time from the double buffer
- Remove deq[i] as an output port: deq[i] will be derived from en and empty_n_w_[i] signals
- accum_out_skew_fifo.v
 - **When en = 1, each row dequeues 1 data until the fifo becomes empty**
 - only 1 en signal for the n-th fifo
- input_skew_fifos.v
 - each row contains a chain of fifos
 - **normally, each cycle, the data moves from the n-th fifo to the n+1 th**
 - 1 pos clk edge after `enqueue[1]` goes high, data appears at `d_w[1][1]` -- the output of the first fifo in that row
 - empty_n and full_n signals are **combinational**, changes within the same cycle
 - `empty_n[row] = empty_n[row][last]` : if the last fifo is NOT empty, we can dequeue from the chained fifo
 - `full_n[row] = full_n[row][0]`: if the first fifo is not full, we can enqueue to it
 - **The fifo output stays unchanged**, w/o active “dequeue” signal



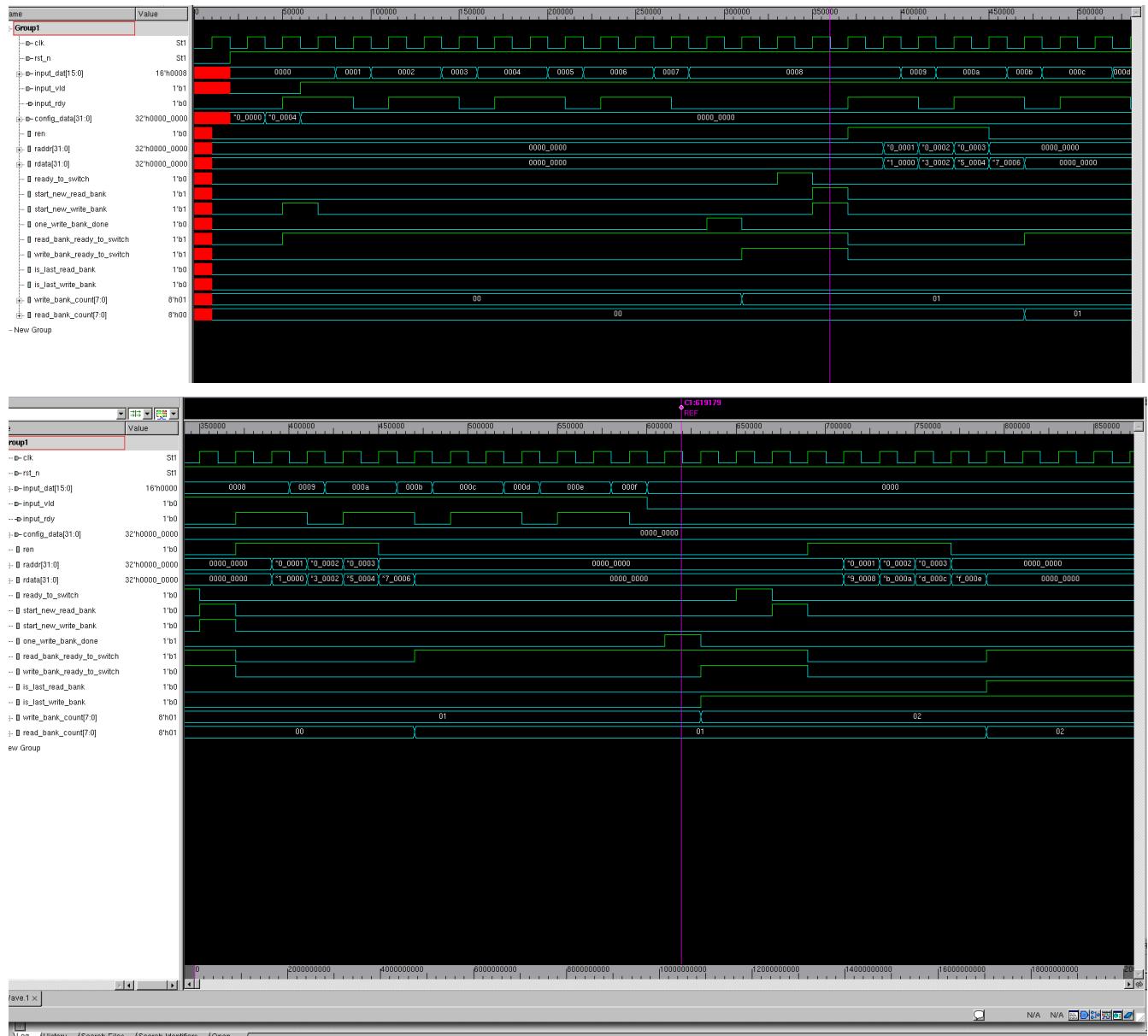
- fifo.v
 - The queue only holds **1D data**.

- What's the COUNTER_WIDTH? The # of elements the queue can save
 - $\log_2(\text{depth}-1)$
 - it creates a fifo module using the given SizedFIFO.

- SizedFIFO.v
 - it's given to us with copyright to Bluespec, Inc.
 - what's p2depth?
 - The number of elements that can be saved in the FIFO.
 - the FIFO array size is actually ($p2\text{depth} - 1$), the one extra element floats
 - p3cntr_width is the binary counter's width to hold the index of the FIFO elements.
 - head points to the element to be dequeued next
 - can hold ($p2\text{depth}-1$) number of integers, it's like an index
 - tail points to the element that was (just enqueued)
 - array saves a number of ($p2\text{depth}-1$) data with p1width bits each.
 - Do we need to set BSV_RESET_FIFO_ARRAY?
 - this clears the array internal values upon rst_n signal.
 - Do we need to set BSV_RESET_FIFO_HEAD?
 - this sets D_OUT to be zero upon reset
 - I think we need this, because rn D_out defaults to aaaa....

Simple main controller

- simple_read_FSM.v
- simple_read_controller.v
 - contains simple_read_FSM, read_addr_gen, read_bank_counter modules
 - relays control signals between the read_FSM and the main_FSM
- simple_main_FSM.v
- simple_main_controller.v
 - write banks of data to the double buffer, and read all banks out in the same order
 - contains ifmap_input_controller, simple_read_controller, and simple_main_FSM.
- simple_main_controller_tb.sv
 - write all banks and then read all banks in the same addr order, for IC0 = 2 (chaining length), IC1*Y0*IX0 = 4 (# of data per bank), OY1*OX1 = 2 (total # of banks)
 - see the simulation waveform below
 - Due to chaining, writing one bank takes longer than reading one.
 - In the CNN, we can read one bank multiple times (i.e. the ifmap data) before switching and writing the next bank.



Testbench files

C files

- conv_gold_sv.cpp
 - the **gold model** that can be called by system verilog modules.

- the callable function is “run_cov_gold()”
 - This file calls “conv_gold.cpp”.
 - **Define CNN parameters** here, needs to match with the SV testbench

- conv_gold.cpp
 - C++ version implementation of the “conv” module
 - is a parameterized function

- conv_gold_test.cpp
 - **Testbench** for “conv_gold.cpp”.
 - It runs 3xx CNN test cases
 - It contains a function “run_layer()” that runs a specific layer using “conv_gold.cpp”, and then compares the results to pre-computed correct values.

Verilog files

- simple_tb_1.sv
 - test bench for **conv_tiled.sv**
 - only 1 ifmap/weight/ofmap bank
 - data in ifmap/weight inputs are positive integers <10.
 - IC0 = OC0 = 2

```

1 simple_tb_1.sv
1 // Description: testbench for conv_tiled.v
2 //           - only 1 ifmap/weight/ofmap bank
3 //           - just directly feeds the ifmap/weight data
4 //           - OY0 = 0X0 = 3, OC0 = IC0 = 2, FY = FX = 2
5 // Author: Cheryl (Yingqiu) Cao
6 // Date: 2022-07-24
7
8 module simple_tb_1;
9
10 // local parameters ++
11 localparam PARAM_NUM = 6;
12 localparam PARAM_WID = 16;
13 localparam BANK_ADDR_WIDTH = 32;      // width needed to save IC1*IX0*IY0, OY1_0X1
14 localparam BUFFER_MEM_DEPTH = 256;    // capacity of the memory, larger than IC1*Ix0*IY0
15
16 // cnn parameters
17 localparam OY0 = 3;
18 localparam 0X0 = 3;
19 localparam OC0 = 2;
20 localparam IC0 = 2;
21 localparam Stride = 1;
22
23 // parameters on: input data dimensions
24 localparam OY = 3;
25 localparam 0X = 3;
26 localparam OC = 2;
27 localparam IC = 2;
28 localparam FY = 2;
29 localparam FX = 2;
30
31 // derived parameters

```

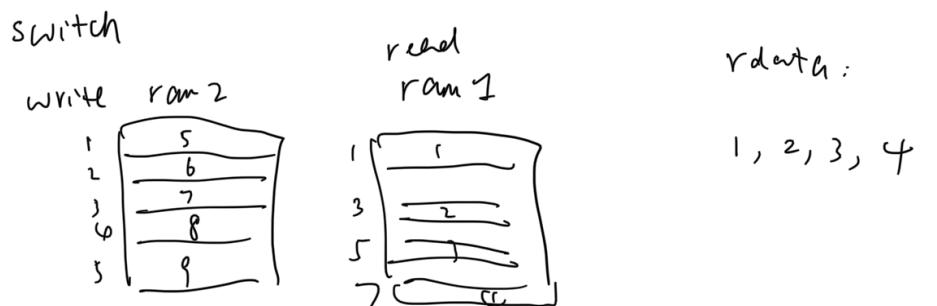
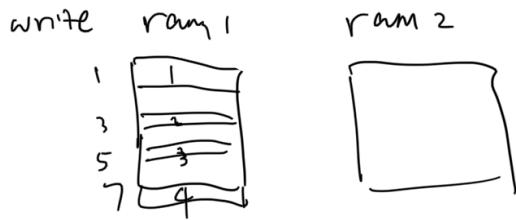
- simple_tb_2.sv
 - test bench for **conv_tiled.sv**
 - similar to tb_1, but contains negative numbers in ifmap/weight inputs
- conv_tb.sv
 - test bench for **conv.sv**
 - conv.sv module computes ofmap data, which will be saved in the transaction item
 - Scorebox calls “run_conv_gold()” to compute **gold_ofmap**, and compared ofmap data with the correct values.
 - The test signal was applied through “test” class => “**apply_sim**” task

```
+1 conv.v | 4 ifmap_input_controller.v | 5 layer_params.svh | 6 conv_tb.sv | 7 conv_gold.cpp | 8 conv_gold_test.cpp | 11 conv.
310
311
312 module conv_tb;
313
314   reg clk;
315
316   always #10 clk =~clk;
317
318   conv_if _if (clk);
319
320   conv #(`IFMAP_SIZE(`LAYER_IFMAP_SIZE), `WEIGHTS_SIZE(`LAYER_WEIGHTS_SIZE), `OFMAP_SIZE(`LAYER_OFMAP_SIZE))
321   dut (
322     .clk(_if.clk),
323     .rst_n(_if.rst_n),
324     .ifmap_dat(_if.ifmap_dat),
325     .ifmap_rdy(_if.ifmap_rdy),
326     .ifmap_vld(_if.ifmap_vld),
327     .weights_dat(_if.weights_dat),
328     .weights_rdy(_if.weights_rdy),
329     .weights_vld(_if.weights_vld),
330     .ofmap_dat(_if.ofmap_dat),
331     .ofmap_rdy(_if.ofmap_rdy),
332     .ofmap_vld(_if.ofmap_vld),
333     .layer_params_dat(_if.layer_params_dat),
334     .layer_params_rdy(_if.layer_params_rdy),
335     .layer_params_vld(_if.layer_params_vld)
336   );
337
338   initial begin
339     test t0;
340
341     clk <= 0;
342     _if.rst_n <= 0;
343     #20 _if.rst_n <= 1;
344   
```

- input_read_addr_gen_tb.v
 - test bench for **input_read_addr_gen.v**
 - sets the blocking parameters for one example and verifies the generated addr by the dut.
- double_buffer_tb.v
 - test read/write functions for the double buffer
 - see the figure below for functionality

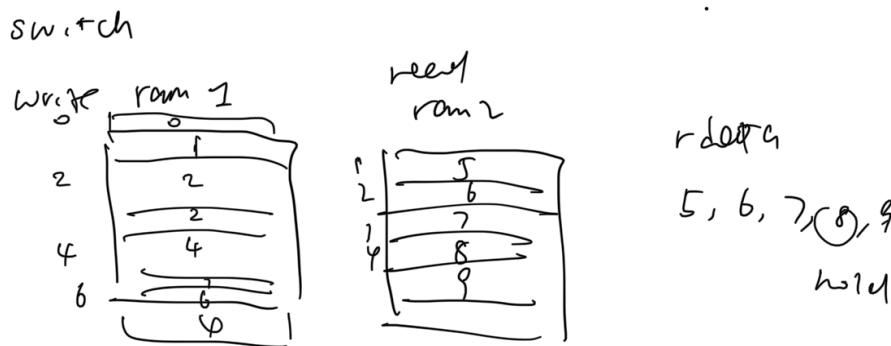
- see the progress notes for 2021-11-28 for the simulation waveforms

double buffer test bench

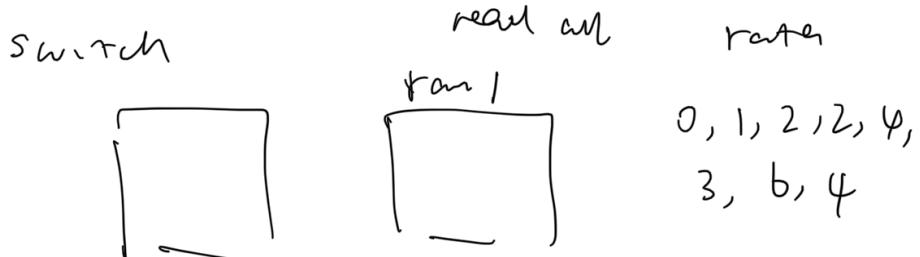


wait 1 cycle

xp

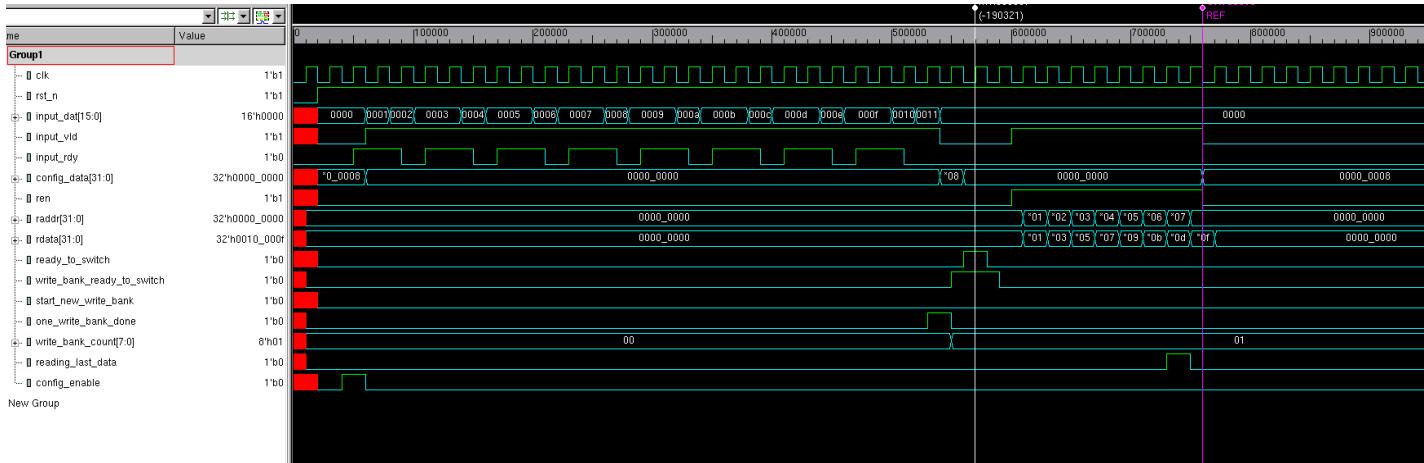


hold

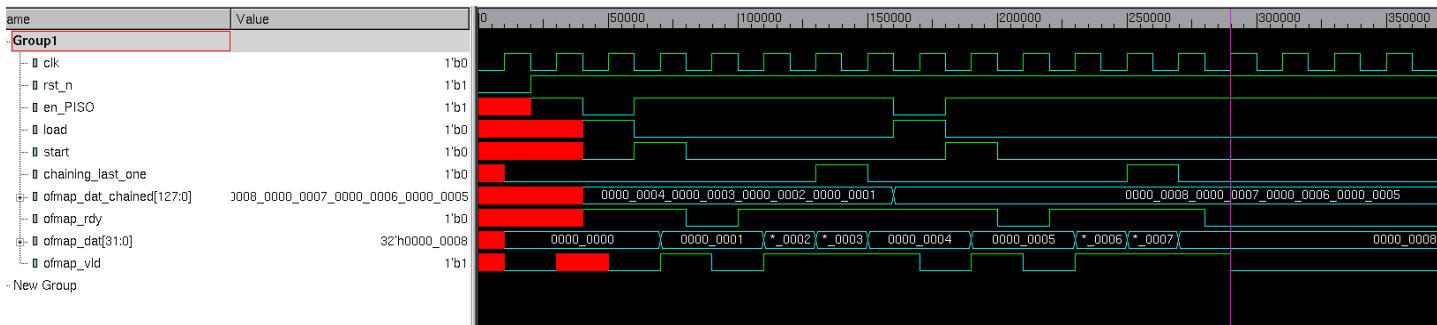


- ifmap_input_controller_tb1.sv
 - Tests the ifmap FSM, write/read address_gen, and double buffer
 - Writes one ifmap bank with IC0 = 2, IC1*IX0*IY0 = 2*2*2 = 8, and then read it for once

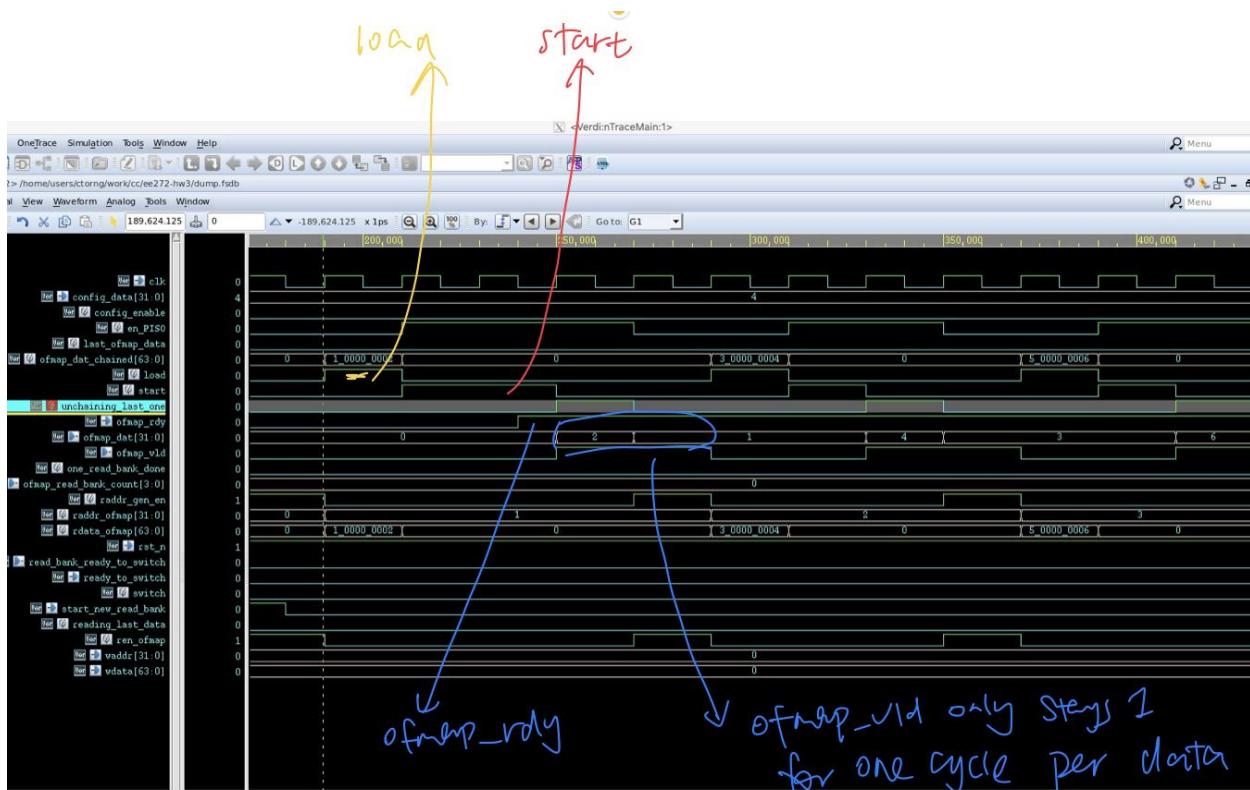
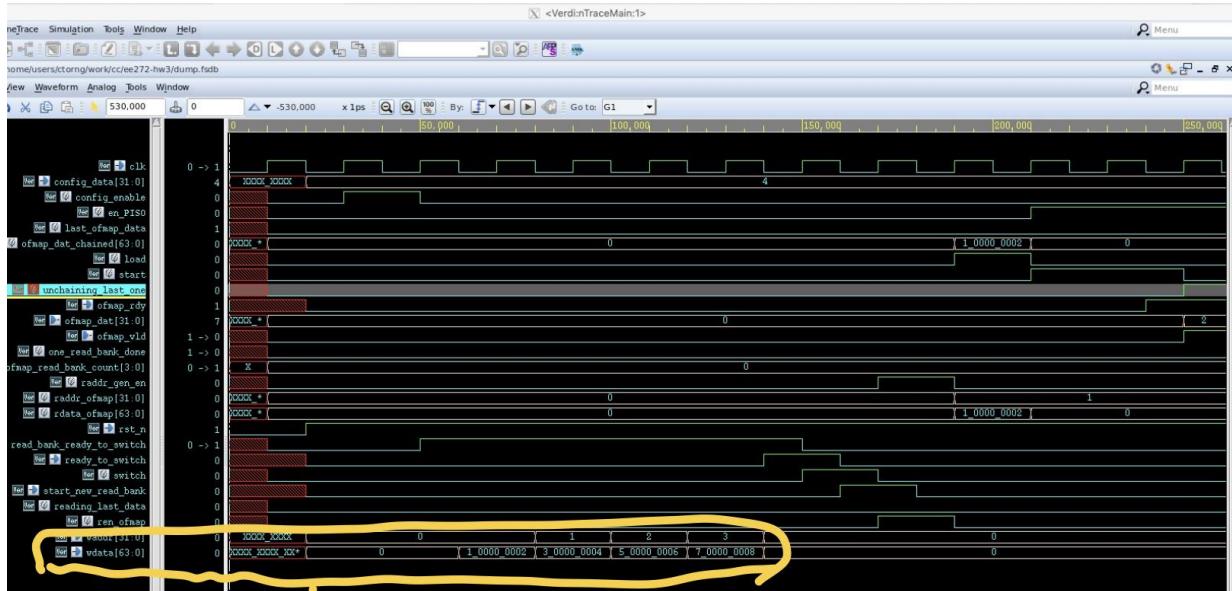
See the simulated waveform as below.

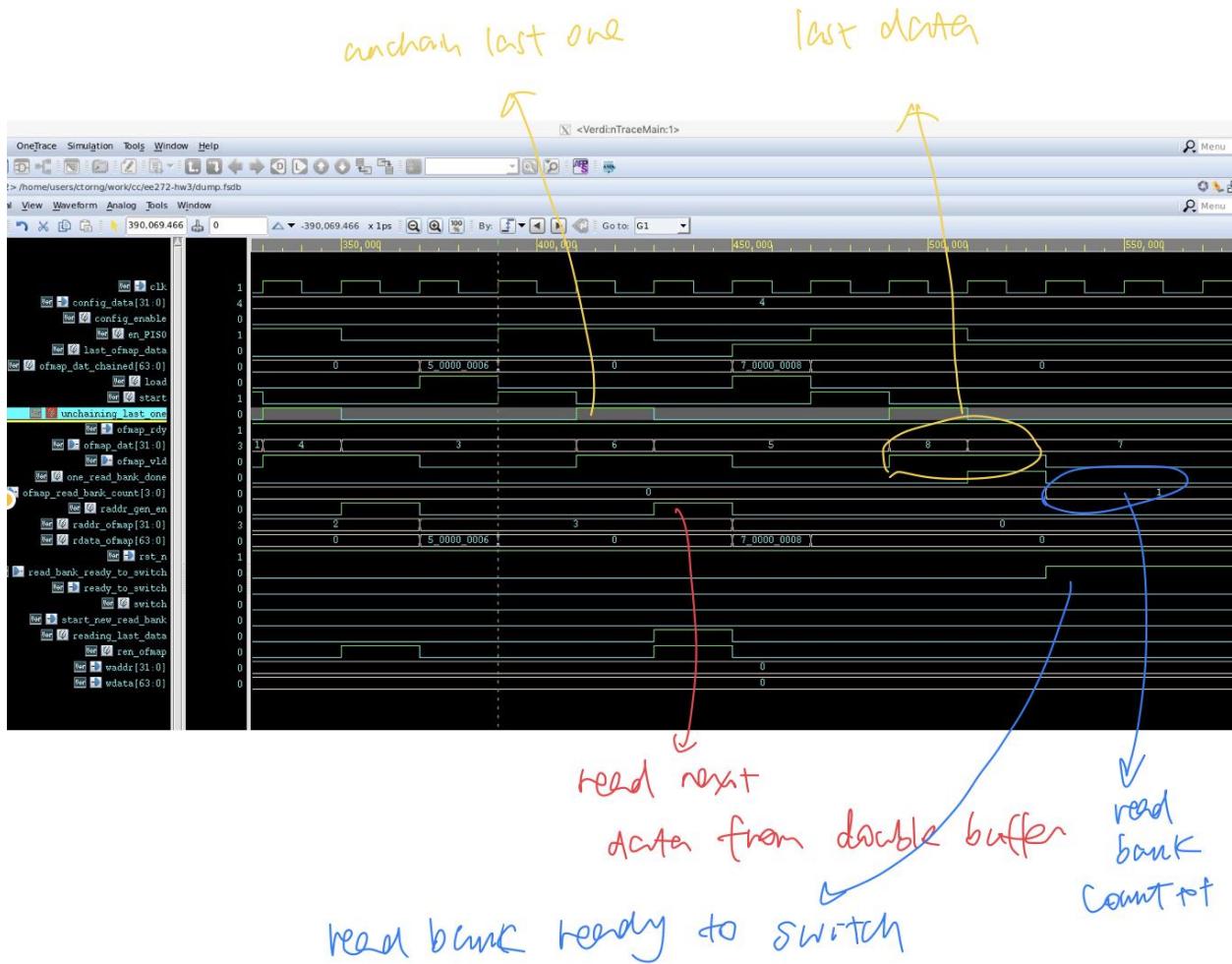


- `ofmap_PISO_tb.sv`
 - OC0 = 4, 2 unchain cycles
 - unchains 4_3_2_1 and 8_7_6_5



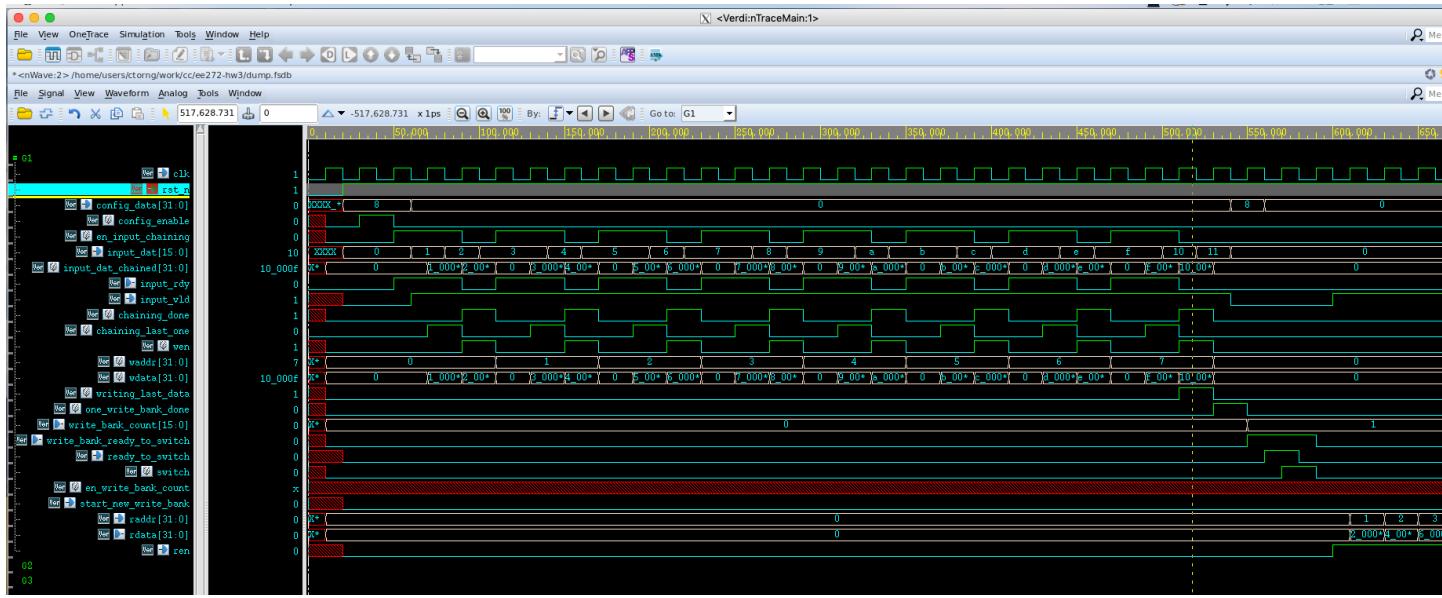
- `ofmap_output_controller_tb1.sv`
 - OY0_OX0 = $2^2 \cdot 2 = 4$
 - OC0 = 2 for unchaining
 - run for one ofmap read bank
 - white data 1_2, 3_4, 5_6, 7_8 to `accum_double_buffer`
 - switch banks
 - let ofmap_output_controller run to read ofmap data out and unchain them
 - monitor ofmap_data, we should get 2,1,4,3, ..., 8, 7.





- weight_input_controller_tb1.sv
 - Writes one weight bank with OC0 = 2, OC1*IC1*Fy*Fx*IC0 = 2*2*2 = 8, and then read it for once

See the simulated waveform as below.



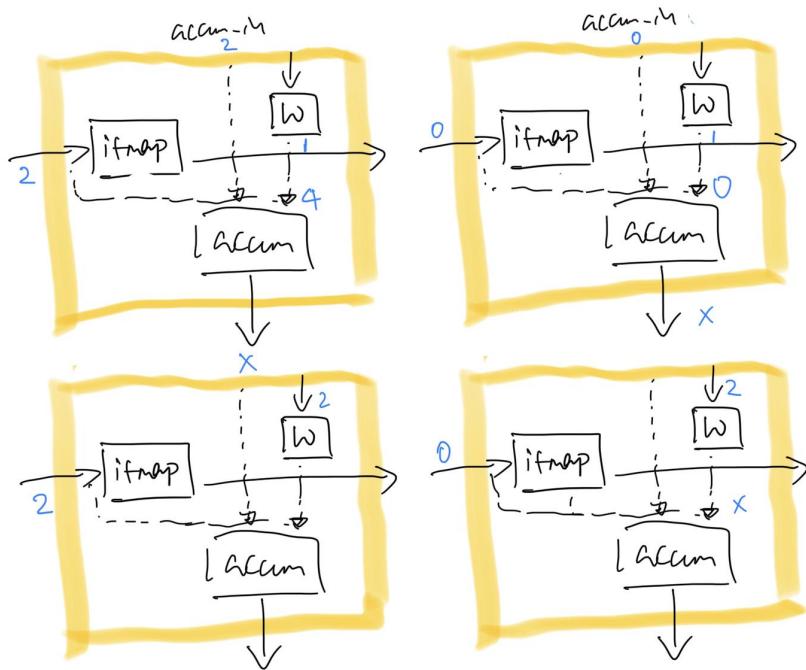
- mac_tb.sv

```
+1 mac_tb.sv | 3 mac.v | 4 ofmap_output_controller_tb1.sv
 1 // Description:
 2 //           ifmap_in   weight_in   accum_in   en_weight   |   accum_out
 3 //           x          1          x          1          |   x
 4 //           1          1          1          0          |   2
 5 //           2          2          2          1          |   4
 6 //           3          3          3          0          |   9
 7 //           4          4          4          1          |   12
 8 //           5          5          5          0          |   25
 9 //       weight_in needs to be ready 1 cycle earlier than ifmap_in and
10 //       accum_in
11 // Author: Cheryl (Yingqiu) Cao
12 // Date: 2022-04-2
13
14
```

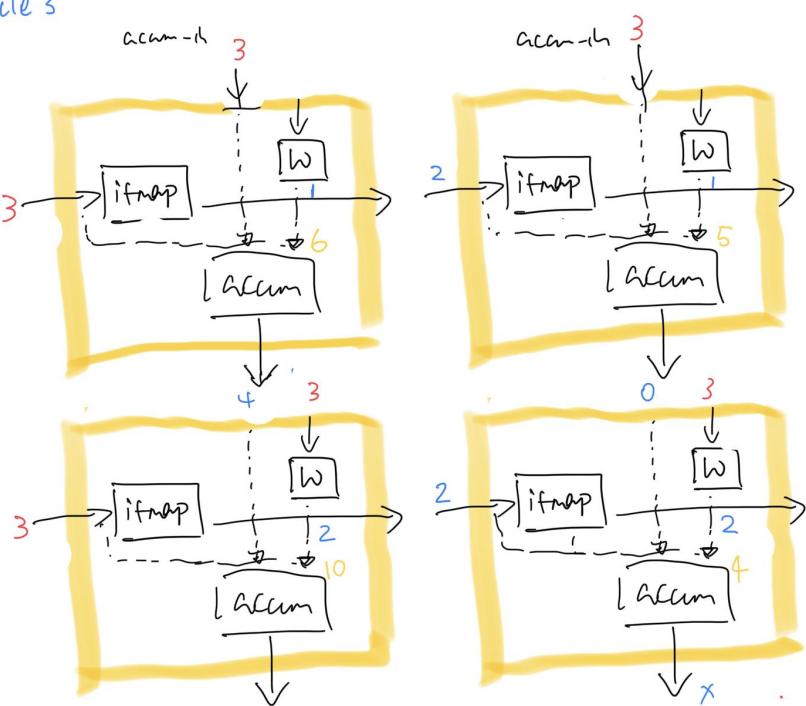
- mac_array_tb.sv

```
1 mac_array_tb.sv | 3 Makefile
1 // Description:
2 // IC0 = OC0 = 4
3 // en_weight goes high for each row one by one (2 cycles) and then turn
4 // off
5 // chaining goes MSB -> LSB
6 // en_weight weight_chained ifmap_chained accum_in | accum_out
7 //      1          1_1          0_0          0_0 | █
8 //      1          2_2          2_2          0_2 |
9 //      0          3_3          3_3          3_3 |
10 //     0          4_4          4_4          4_4 | 4_10
11 //     0          0            0            0 | 11_14
12 //     cycle 6                                | 15_8
13 //     cycle 7                                | 4_0
14 //     cycle 8                                | 0_4
15 // Author: Cheryl (Yingqiu) Cao
16 // Date: 2022-04-03
17
18
```

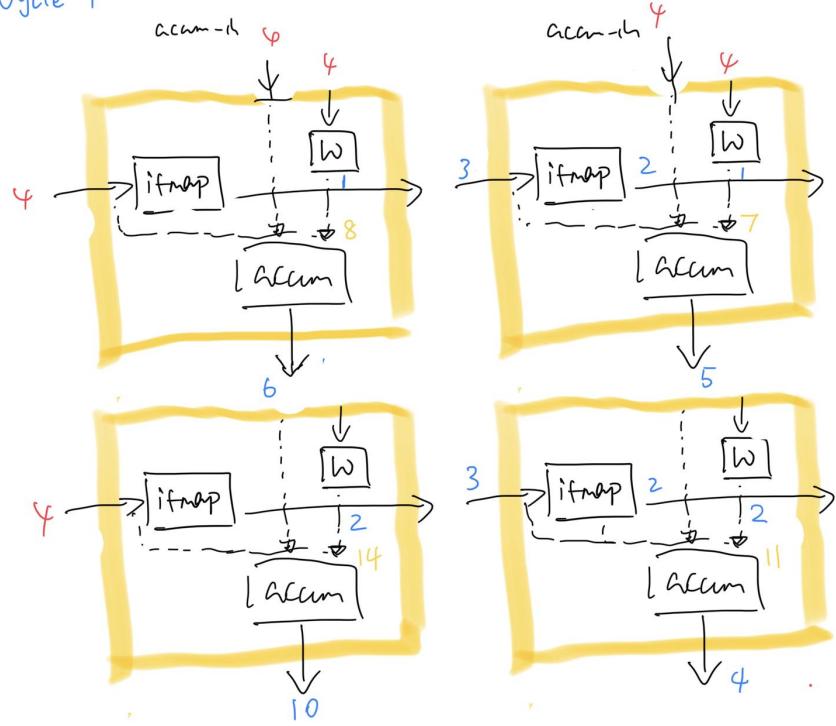
Cycle 2



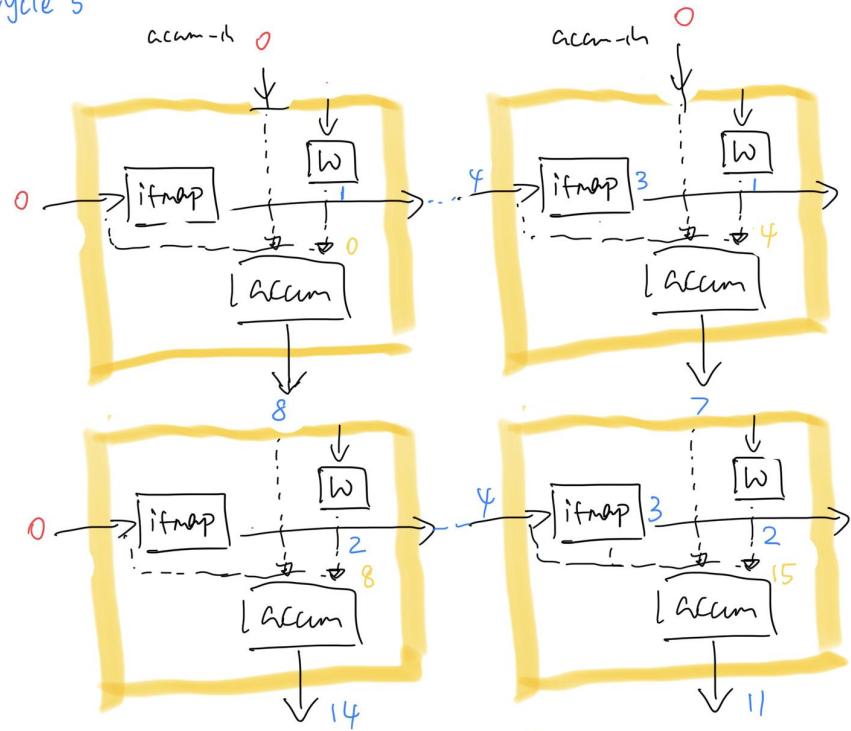
Cycle 3



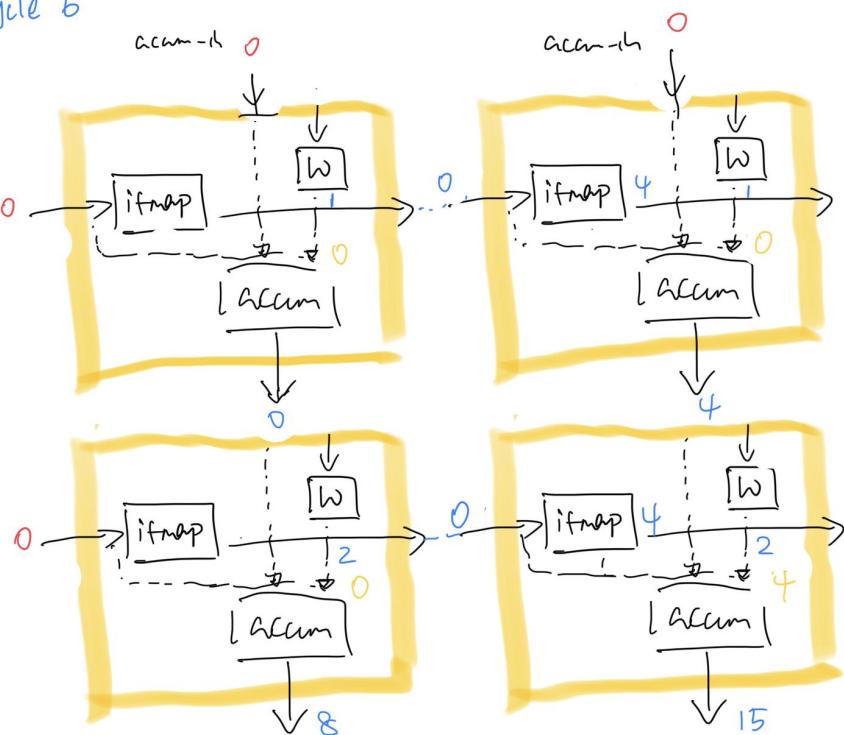
Cycle 4

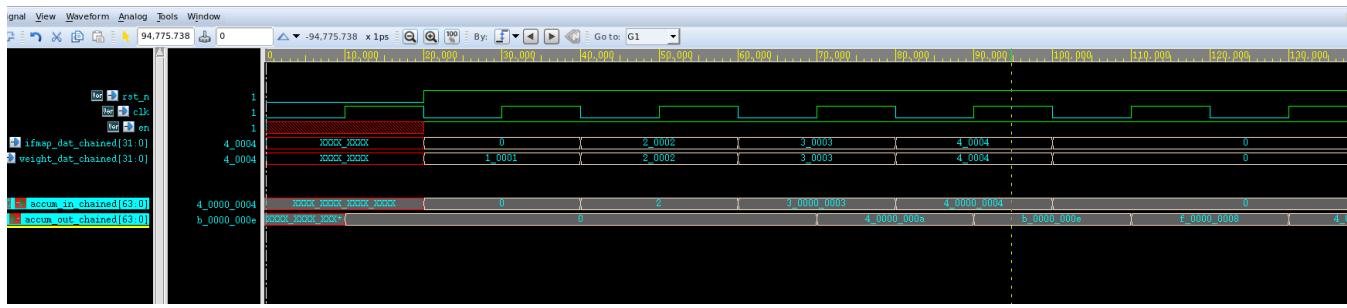


Cycle 5



Cycle 6





Typedef files

- layer_params.svh
 - defines user type “layer_params_t”

Verilog Notes

- Need to use non-blocking assignment “`<=`” inside of always blocks
 - `<=` means to evaluate the expression AT the clock edge
 - the exception is `always(*)`, which is essentially combinational, in which case it does not matter.
- “Casez” allows you to check for multiple bits at the same time:

```
casez ({CLR, DEQ, ENQ, hasodata, ring_empty})
  // Clear operation
  5'b1????: begin
    head      <= `BSV_ASSIGNMENT_DELAY {p3cntr_width{1'b0}} ;
    tail      <= `BSV_ASSIGNMENT_DELAY {p3cntr_width{1'b0}} ;
    ring_empty <= `BSV_ASSIGNMENT_DELAY 1'b1;
    not_ring_full <= `BSV_ASSIGNMENT_DELAY 1'b1;
    hasodata   <= `BSV_ASSIGNMENT_DELAY 1'b0;
  end
  // -----
  // DEQ && ENQ case -- change head and tail if added to ring
  5'b011?0: begin
    tail      <= `BSV_ASSIGNMENT_DELAY next_tail;
    head      <= `BSV_ASSIGNMENT_DELAY next_head;
  end
  
```

- Head `<= delay next_head` assigns the signal with delay

It's equivalent to `head <= #5 next_head`

- `{}` is for concatenation. The following means to repeat the inner bracket for p3cntr number of times

```
assign next_head = (head == depthLess2) ? {p3cntr_width{1'b0}} : incr_head ;
assign next_tail = (tail == depthLess2) ? {p3cntr_width{1'b0}} : incr_tail ;
```

- Macro assertion
 - here if the macro `BSV_RESET...` is not set, the compiler will ignore the lines between `'ifdef` and `'endif`

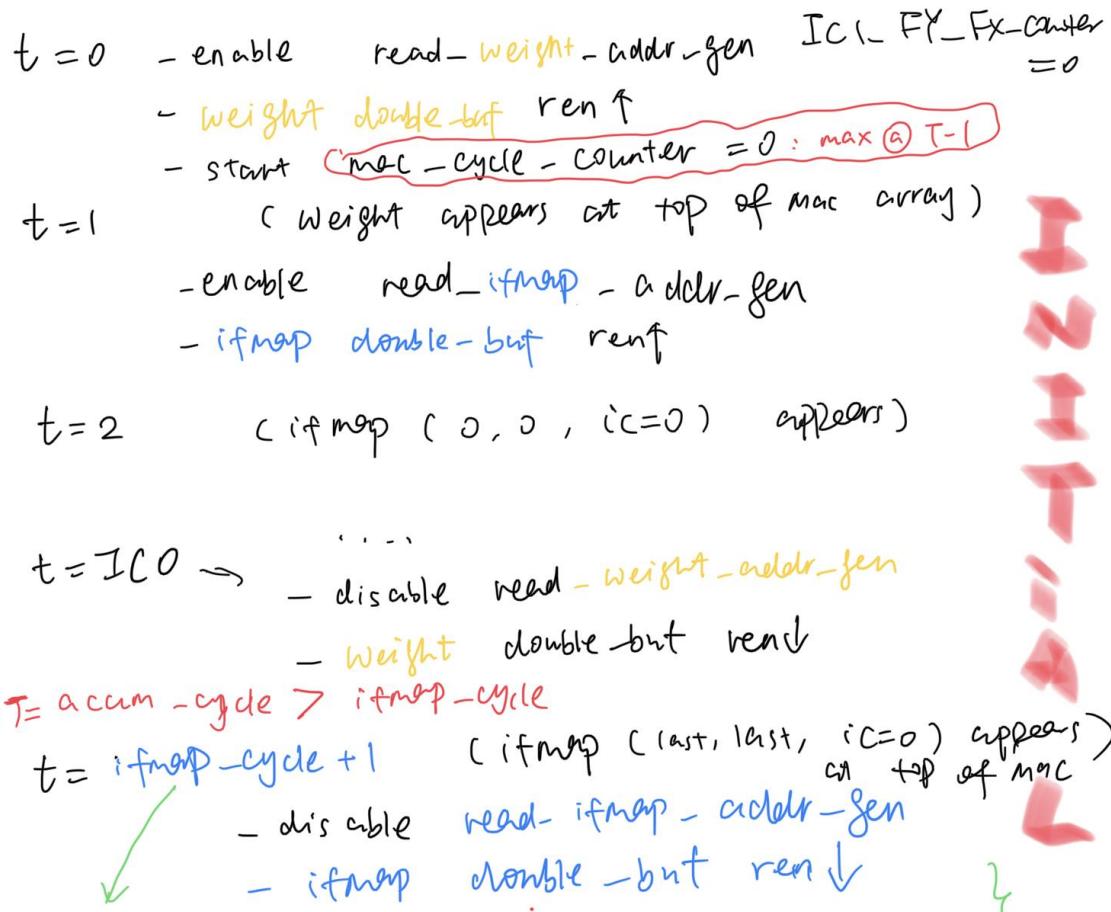
```

// Update the fast data out register
always @(posedge CLK `BSV_RESET_EDGE_HEAD)
begin
`ifdef BSV_RESET_FIFO_HEAD
if (RST == `BSV_RESET_VALUE)
begin
| D_OUT <= `BSV_ASSIGNMENT_DELAY {plwidth {1'b0}} ;
end // if (RST == `BSV_RESET_VALUE)
else
`endif
begin
casez ({CLR, DEQ, ENQ, hasodata, ring_empty})
// DEQ && ENQ cases
5'b011?0: begin D_OUT <= `BSV_ASSIGNMENT_DELAY arr[head]; end
5'b011?1: begin D_OUT <= `BSV_ASSIGNMENT_DELAY D_IN; end
// DEQ only and data is in ring

```

Appendix

- Synchronization case 1: accum_cycle > ifmap_cycle
 - It takes longer for the 1st partial sum ofmap(0,0, IC = 0:3) to get out of the MAC array than to iterate through the ifmap elements for weight (Fx = 0, Fy = 0)
 - need to pause ifmap double buffer input



$t = T+2$ (partial sum of map $(0,0)$, $iC=0:3$) appears (5)
 top of mac array.
 ifmap $(1,0, iC=0)$ - for $Fx=1, FY=0$ appears)

$t = ICO$ same as Initial

$t = \text{ifmap-cycle}$ (partial sum of map $(\text{left}, \text{out}, iC=0:3)$
 ready to de queue)
 - disable write-accum-addr-gen
 - accum-buf wen ↓
 - plus operations in Initial

Normal

$t = \text{ifmap-cycle} + 1$
 - disable read-accum-addr-gen
 - accum-buf wen ↓
 - plus operations in Initial

Final

$$\begin{aligned} & \text{IC1-FY-FX counter} \\ & = = \text{IC1} \times \text{FY} \times \text{FX} \\ & 2 \times 3 \times 3 = 18 \end{aligned}$$

$t = 0$ instead of the accum-buffer
 write to the ofmap-double-buf
 - enable write-accum-addr-gen
 - ofmap-double-buf wen ↑

$t = \text{ifmap-cycle}$ (ofmap write complete)

- disable write-accum-addr-gen
 - ofmap-double-buf wen ↓

one_ifmap_bank_done

- reset mac array
- switch ifmap / ofmap double-buf
(wait if ifmap write / ofmap read not complete yet)
- increment ifmap-bank counter
 $c++ @ \text{next CLK cycle}$
max @ $0Y1 \times 0X1 \times 0C1$
 $3 \times 3 \times 3 = 27$

- Synchronization case 2: ifmap_cycle > accum_cycle
 - It takes longer to iterate through the ifmap elements for weight ($F_x = 0, F_y = 0$), than for the 1st partial sum ofmap($0,0, IC = 0:3$) to get out of the MAC array
 - need to wait to feed back the accum_buf to MAC array

- (2) $T = \text{ifmap-cycle} > \text{accum-cycle}$
- $\text{o}_Y_0 \cdot \text{o}_X_0 = 9$ " $\text{o}_{CO} + \text{I}_{CO} + 1 = 9$
 $2 \cdot 2 = 4$ " $2 + 2 + 1 = 5$
- $t=0$ - enable read-weight-addr-gen $I_{CL-FK-Fx-Counter} = 0$
 - weight double-buf ren↑
 - start Mac-cycle-counter = 0: max @ $T-1$
 $t=1$ (Weight appears at top of Mac array)
 - enable read-ifmap - addr-gen
 - ifmap double-buf ren↑
 $t=2$ (ifmap(0,0,ic=0) appears)
- INITIAL
- $t=I_{CO} \rightarrow$ - disable read-weight-addr-gen
 - weight double-buf ren↓
- $t=\text{accum-buf-cycle}$ (1st partial sum of ifmap(0,0,ic=0:3)
 ready to enqueue)
 - enable write-accum-addr-gen
 - accum-buf wen↑
 - read weight ($FX=1, FK=0$)
- $t=\text{ifmap-cycle}/0$
 - enable read-weight-addr-gen
 - weight double-buf ren↑
 - Mac-cycle-Counter rolls over to 0
 - increment I_{CL-FK-Fx-Counter}
- $t=\text{ifmap-cycle} + 1$ (ifmap(last, last, ic=2) appears
 @ top of Mac array)
 - enable read-accum-addr-gen
 - accum-buf ren↑
- DECODE

- $t=2$
- ifmap(0, 0, i[=0]) for $Fx=1, Fy=0$ appears
 - ofmap(0, 0, i[=0]) @ top of MAC
 - weight $Fx=1, Fy=0$ ready)

$t=7co$

same as (Initial)

Normal

$t = \text{accum-but-cycle}$

same as

Initial

Pre-final

- Other cycles same as (normal)

$t = \text{accum-but-cycle}$ (1st partial sum of map(0, 0, ..))
ready to enqueue

write to ofmap-double-buf

- enable write-accum-addr-gen
- ofmap-double-buf wen↑

Final

ICL-FY-FX counter

$$= = ICL \times FY \times FX$$

$$2 \times 3 \times 3 = 18$$

$t=0$ stop reading weights

- disable read-weight-addr-gen

- weight-buf ren↓

$t=1$ stop reading ifmap (ifmap(lwt), lwt) appears
@ top of MAC

- disable read-ifmap-addr-gen

- ifmap-buf ren↓

$t = \text{accum_buf_cycle} - 1$ (ofmap((last), last) ready
to enqueue)

$t = \text{accum_buf_cycle}$ (ofmap - writing done)

- disable write - accum-adder-gen

- ofmap - double-buf wen↓

one-ifmap-bank-done

- reset MAC array

- switch ifmap / ofmap double-buf
(wait if ifmap write / ofmap read not
complete yet)

- increment ifmap-bank Counter
(++ at next CLK cycle)

$$\begin{array}{rcl} \text{max @ } & 0Y1 \times 0X1 \times 0C1 \\ & 3 \times 3 \times 3 = 27 \end{array}$$

- Synchronization case 3: ifmap_cycle == accum_cycle
 - It takes the same # of cycles to iterate through the ifmap elements for weight ($F_x = 0, F_y = 0$), as for the 1st partial sum ofmap(0,0, IC = 0:3) to get out of the MAC array.
 - A special case of Case 2.

$$③ T = \text{ifmap-cycle} = \text{accum-cycle}$$

$$\text{ifmap-cycle} = 9$$

$$OY_0 \cdot OX_0 = 9$$

$$OC_0 + IC_0 + 1 = 9$$

$$2 + 2 + 1 = 5$$

$$IC_1 - F_X - F_X - \text{center} = 0$$

- $t=0$
- enable read-weight-addr-gen
 - weight double-buf ren↑
 - start Mac-cycle-Counter = 0 : max @ $T-1$
- $t=1$
- (Weight appears at top of Mac array)
 - enable read-ifmap - addr-gen
 - ifmap double-buf ren↑
- $t=2$
- (ifmap (0, 0, iC=0) appears)

INITIAL

$$t=IC_0 \rightarrow \dots$$

- disable read-weight-addr-gen

- weight double-buf ren↓

-
- $t=$ ~~IC₀~~ ^T accum-buf-cycle $\text{ifmap-cycle}/0$ (1st partial sum of ifmap(0, 0, iC=0:3))
- ready to enqueue
- enable write-accum-addr-gen
 - accum-buf wen↑

(read weight $F_X=1, F_Y=0$)

- mac-cycle-Counter rolls

over to 0

- increment $IC_1 - F_X - F_X$ Counter

- $t=$ ~~ifmap-cycle~~ $+ 1$ (ifmap (last, last, iC=0) appears @ top of Mac array)
- (degree of ifmap(0, 0, iC=0:3))
- enable read-accum-addr-gen
 - accum-buf ren↑

- $t=2$
- ifmap(0, 0, i=0) for $Fx=1, Fy=0$ appears
 - ofmap(0, 0, i=0, ..) (top of MAC)
 - weight $Fx=1, Fy=0$ ready)

$t=7_{CO}$ same as (Initial) **Normal**

Final

$$IC_1 - FY - FX \text{ counter} \\ = = IC_1 \times FY \times FX \\ 2 \times 3 \times 3 = 18$$

$t=$ accum-buf-cycle (1st partial sum ofmap(0, 0, ..))
ready to enqueue)
write to ofmap-double-buf instead

- enable write-accum-addr-gen
- ofmap-double-buf wen↑
- accum-buf wen↓

- △ stop reading weights
- disable read-weight-addr-gen
- weight-buf ren↓

$t=1$ △ stop reading ifmap (ifmap((out), (out)) appears
(top of MAC))
- disable read-ifmap-addr-gen
- ifmap-buf ren↓
△ stop feeding accum to MAC
- disable read-accum-addr-gen
- accum-buf ren↓

$t=$ accum-buf-cycle -1 (ofmap((last), (last)) ready
to enqueue)

one_ifmap-bank-done

$t = \text{accum-buf-cycle}^{10}$ (ofmap-writing done)
— disable write-accum-adder-gen
— ofmap-double-buf wen↓

- reset Mac array
- switch ifmap / ofmap double-buf
(wait if ifmap write / ofmap read not complete yet)
- increment ifmap-bank counter
(++ @ next CLK cycle)
 $\max @ 0Y1 \times 0X1 \times OC1$
 $3 \times 3 \times 3 = 27$