

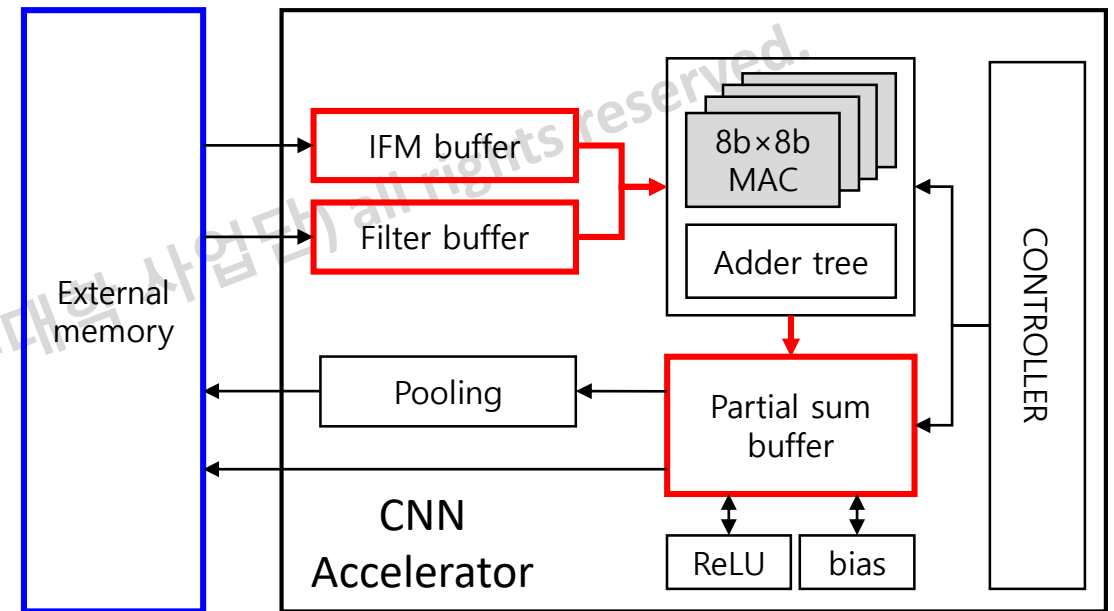
On-chip buffers, Block RAM

2024.03.15 (Thursday)



Motivation

- On-chip buffers (memory)
 - Input feature map (IFM)
 - Filters
 - Intermediate results or partial sum
- External memory



Example of filter, bias, scale files

- Three types of buffers
 - Filter/weight: 16 lines, each line has 128 bits.
 - Bias/scale: 16 lines, each line has 16 bits

conv_weight_L1.hex

	conv_weights_L1.hex	cnn_accel.v	spram.v
1	00000000000000071FFE5CDA37FD7978E		
2	000000000000000879DD5F8013D1B545		
3	0000000000000001EDFCF07FF1FFF40D		
4	000000000000000D334305A2080EB8745		
5	00000000000000033EF09D5FB8003CFF3		
6	000000000000000FEEEE02FC7FF6FAD2F8		
7	000000000000000F880F7165702F4020D		
8	000000000000000FB41CD8BB82B3480F3		
9	00000000000000080AB00093CD20ED39F		
10	000000000000000E2EE10AC4403C180CE		
11	000000000000000FB0501FEBDF2F67F0D		
12	00000000000000060EF6A77FF201FDFE		
13	0000000000000001F815BE67FD4F8800D		
14	000000000000000193EDB928BAE1F3080		
15	0000000000000008BAB80F90D42F02400		
16	0000000000000000AF906E67FA8ED60C6		

conv_bias_L1.hex

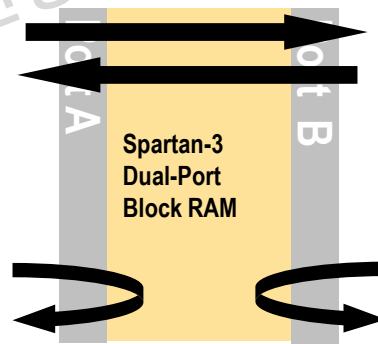
	conv_biases_L1.hex
1	b744
2	1f82
3	0172
4	fe06
5	f911
6	f942
7	171d
8	f98c
9	066a
10	0199
11	f8b5
12	fd24
13	d332
14	e45d
15	0232
16	fe6b

conv_scale_L1.hex

	conv_scales_L1.hex
1	005f
2	0067
3	016c
4	00aa
5	007a
6	0136
7	00cb
8	0079
9	006b
10	00a0
11	0135
12	0107
13	004d
14	006a
15	00af
16	00ba

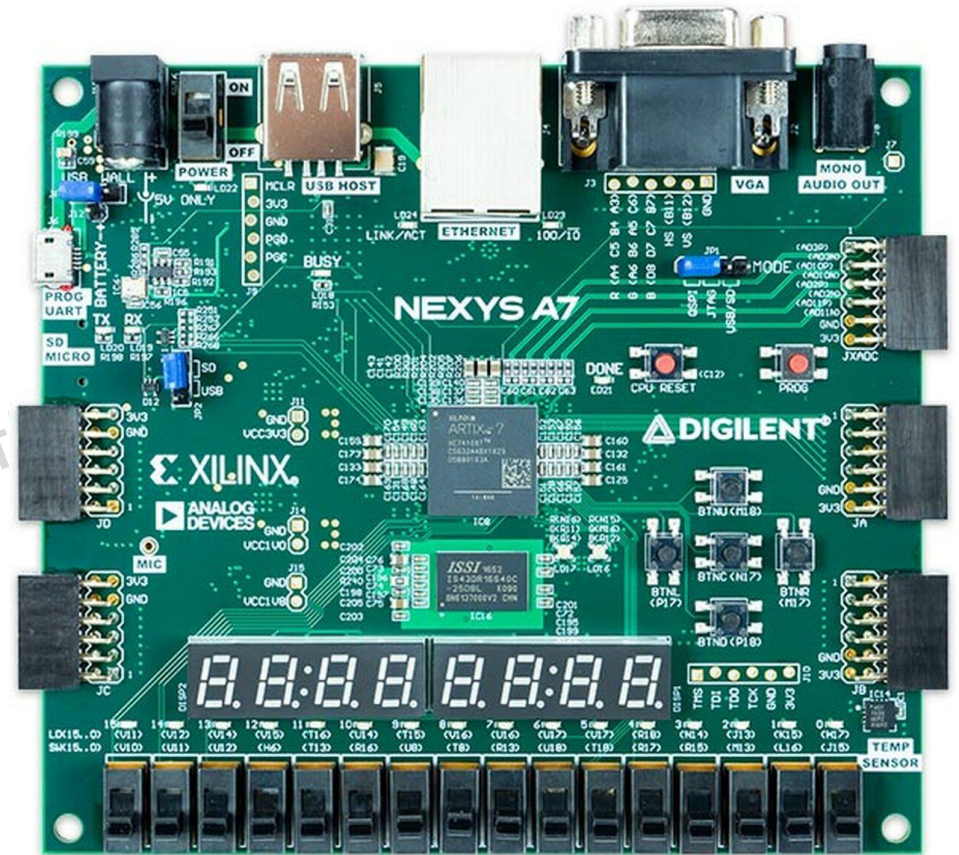
Block RAM

- Most efficient memory implementation
 - Dedicated blocks of memory
- Ideal for most memory requirements
 - 4 to 104 memory blocks
 - **18 kbits = 18,432 bits per block (16 k without parity bits)**
 - Use multiple blocks for larger memories
- Builds both **single** and **true dual-port** RAMs
- Synchronous write and read (different from distributed RAM)



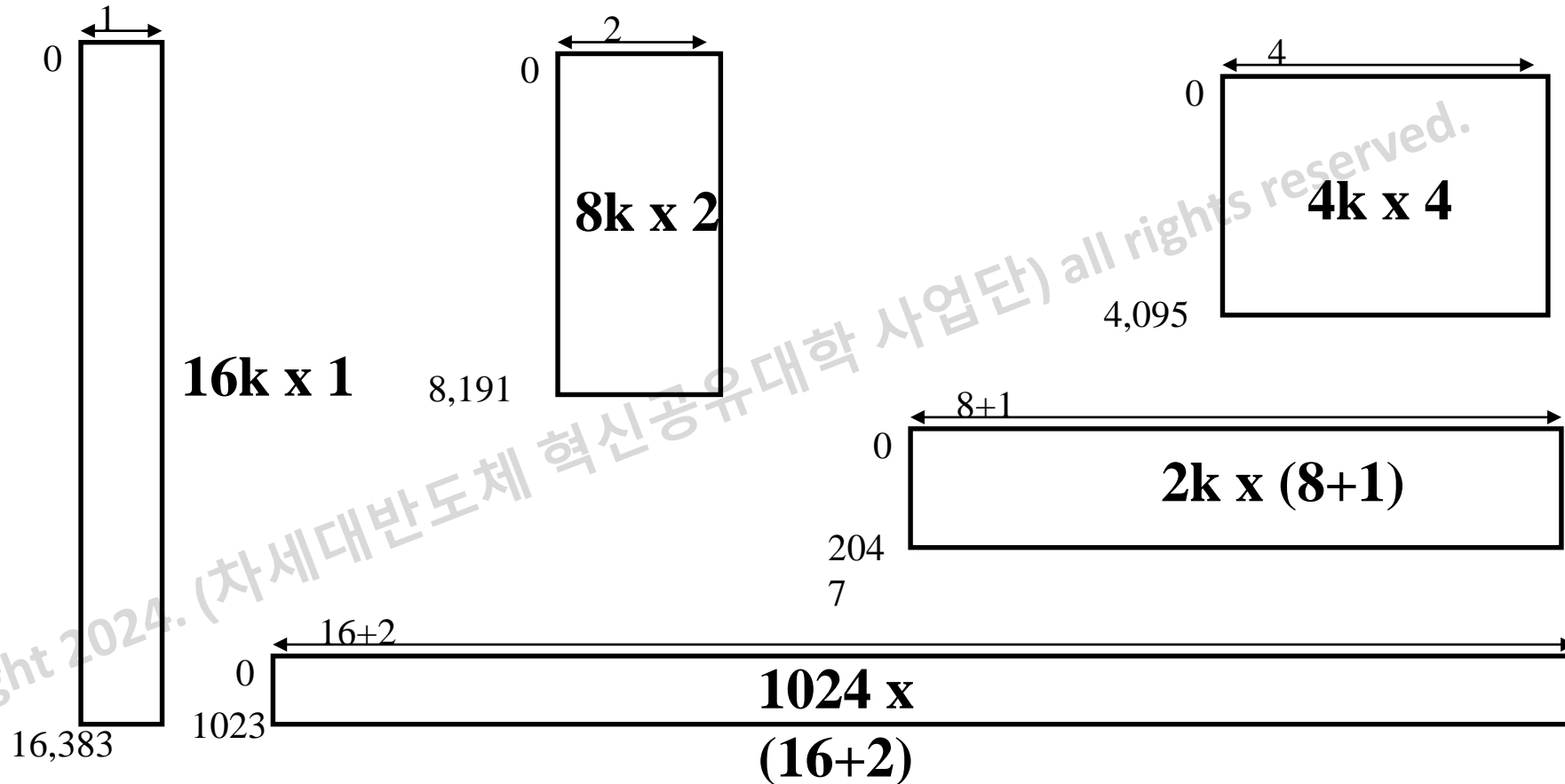
Nexys A7 FPGA board

- Xilinx Artix-7 FPGA XC7A100T-1CSG324C
- 15,850 logic slices
 - Each with four 6-input LUTs and 8 FFs
- 4,860 Kbits of fast block RAM
 - 270 block RAMs
- 240 DSP slices (constrained to To, Ti)
 - Dedicated to multiplication and accumulation (MAC)
- Internal clock speeds exceeding 450 MHz
- 128 MB DDR2 Memory
- USB-JTAG port for FPGA programming and communication



Block RAM Port Aspect Ratios

- Block RAM can have various configurations (port aspect ratios)



Block RAM Port Aspect Ratios

- 18Kb BRAM can be configured
 - RAM_512×36 (512 words, each word has 36 bits).
 - RAM_1K×18 (1024 words, each word has 18 bits).
 - RAM_2K×9 (2048 words, each word has 9 bits).
- 16Kb BRAM can be configured
 - RAM_4K×4 (4096 words, each word has 4 bits).
 - RAM_8K×2 (8192 words, each word has 2 bits).
 - RAM_16K×1 (16384 words, each word has 1 bits).

Organization	Memory Depth	Data Width	Parity Width	DI/DO	DIP/DOP	ADDR	Single-Port Primitive	Total RAM Kbits
512x36	512	32	4	(31:0)	(3:0)	(8:0)	RAMB16_S36	18K
1Kx18	1024	16	2	(15:0)	(1:0)	(9:0)	RAMB16_S18	18K
2Kx9	2048	8	1	(7:0)	(0:0)	(10:0)	RAMB16_S9	18K
4Kx4	4096	4	-	(3:0)	-	(11:0)	RAMB16_S4	16K
8Kx2	8192	2	-	(1:0)	-	(12:0)	RAMB16_S2	16K
16Kx1	16384	1	-	(0:0)	-	(13:0)	RAMB16_S1	16K

Lab 1: Single-port RAM

- Lab 1: Single-port RAM (spram): 16 bits per word, 6240 words
 - Use a single-port RAM
 - Generate a dedicated RAM using Vivado IP generator
 - Initialize a RAM using a coe file
 - Test bench

```
CONV20_W.coe
1 memory_initialization_radix=16;
2 memory_initialization_vector=
3 fe00fe01 ,
4 fd01ffff ,
5 00000303 ,
6 00fdfe02 ,
7 0005ff00 ,
8 00fefd00 ,
9 0000fe00 ,
10 01000000 ,
11 fdff00fb ,
12 01feffff ,
13 ffff0201 ,
14 00fa0001 ,
15 010400fd ,
16 0500ff00 ,
```


Spram and spram wrapper

- Ports

- Clock (clka)
- Read/Write enable (ena)
- Write enable (we)
- Address (addra)
- Write data (dina)
- Read data (douta)

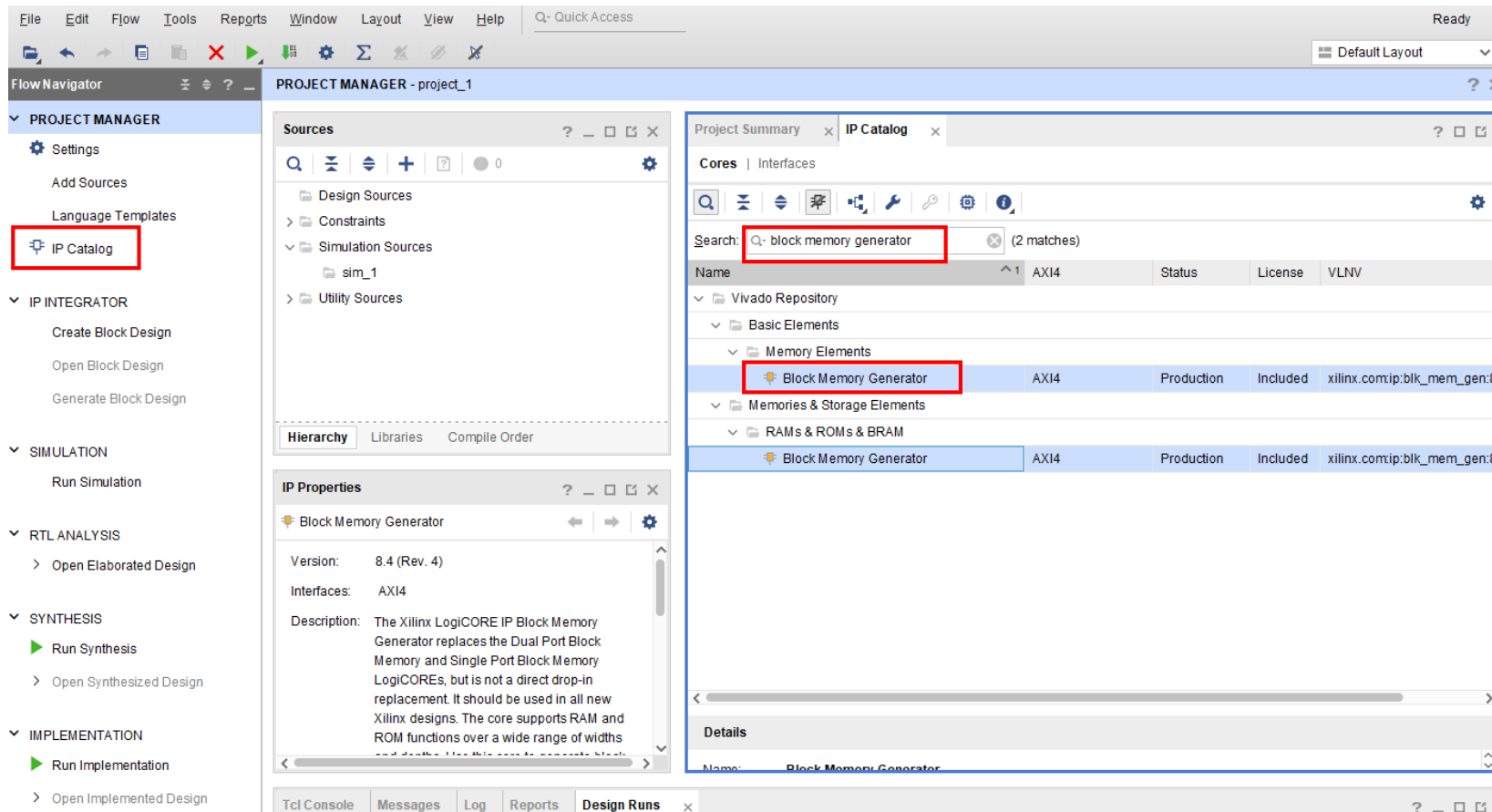
- Two modes

- Simulation
- FPGA
 - Call a SRAM instance generated by IP generator (wrapper)
 - ***Optional: Initialize memory cell from a file***

SPRAM: Simulation

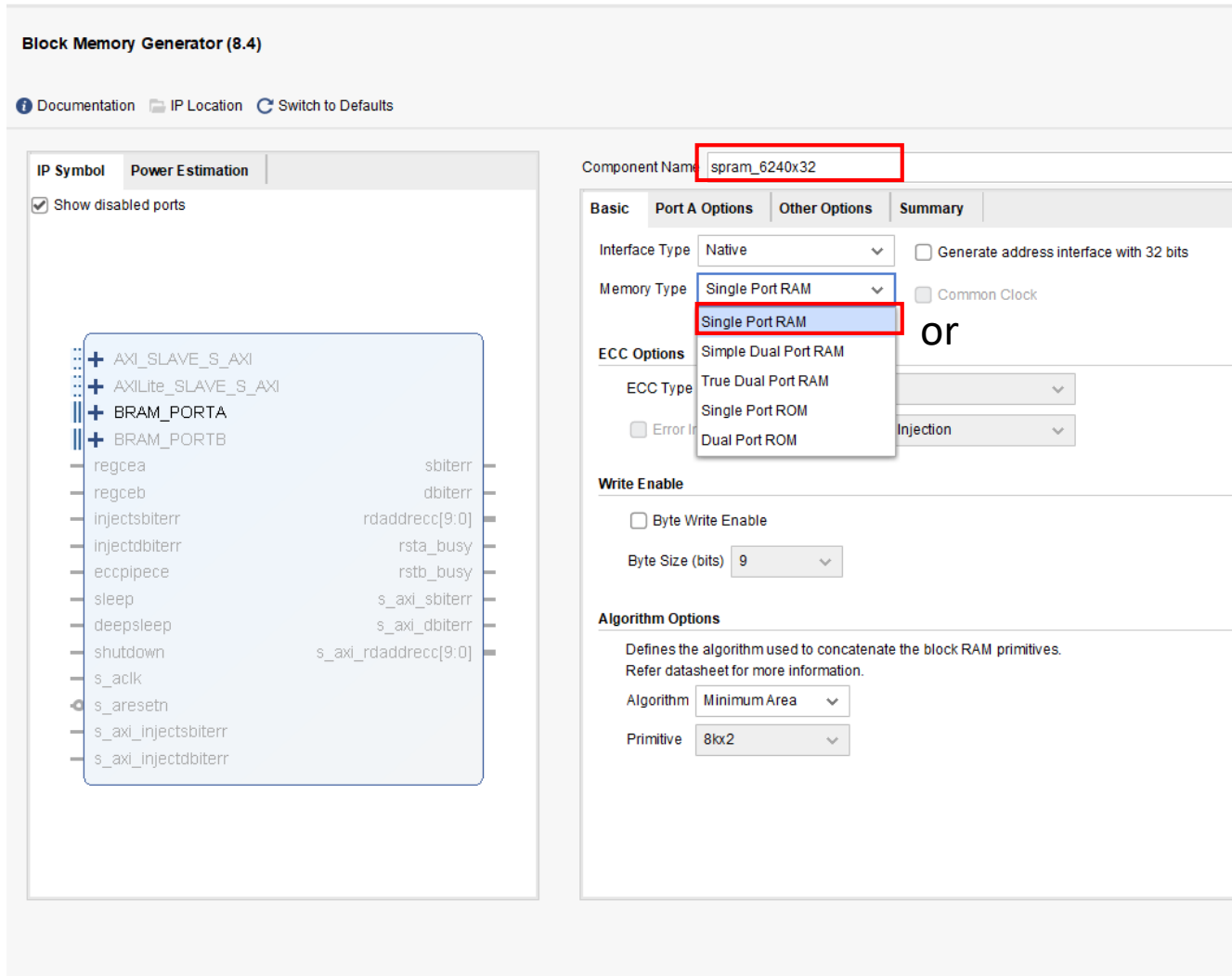
- Memory cell (mem)
 - Data width (DW)
 - Filter: 128, scale/bias: 16
 - Depth
 - Example: 16 lines
- Ports
 - Clock (clka)
 - Read/Write enable (ena)
 - Write enable (wea)
 - Address (addra)
 - Write data (dina)
 - Read data (douta)
- Read operation
 - If(ena)
 - $douta \leq mem[addra]$
- Write operation
 - If(ena & wea)
 - $mem[addra] \leq dina$
- * Single-port
 - Read/write operations share "addra"
 - Only read or write occurs at a time

FPGA: How to make BRAM IP?



- 1. click IP catalog
- 2. search the "block memory generator"
- 3. double click the "block memory generator"

How to make BRAM IP?



- 4. Enter the same name that you declare on your code

- 5. Select the memory type
For reference, in given file, you will use Single Port Ram (spram) and Dual Port Ram(dpram)

How to make BRAM IP?

Component Name

Basic | Port A Options | Other Options | Summary

Memory Size

Write Width Range: 1 to 4608 (bits)

Read Width

Write Depth Range: 2 to 1048576

Read Depth

Operating Mode Enable Port Type

Port A Optional Output Registers

☒ Primitives Output Register ☐ Core Output Register

☐ SoftECC Input Register ☐ REGCEA Pin

Port A Output Reset Options

☐ RSTA Pin (set/reset pin) Output Reset Value (Hex)

☐ Reset Memory Latch Reset Priority

READ Address Change A

☐ Read Address Change A

- 6. Enter the proper width and depth of bram
- 7. Operating mode : No change
- 8. Primitives output register : deselect

How to make BRAM IP?

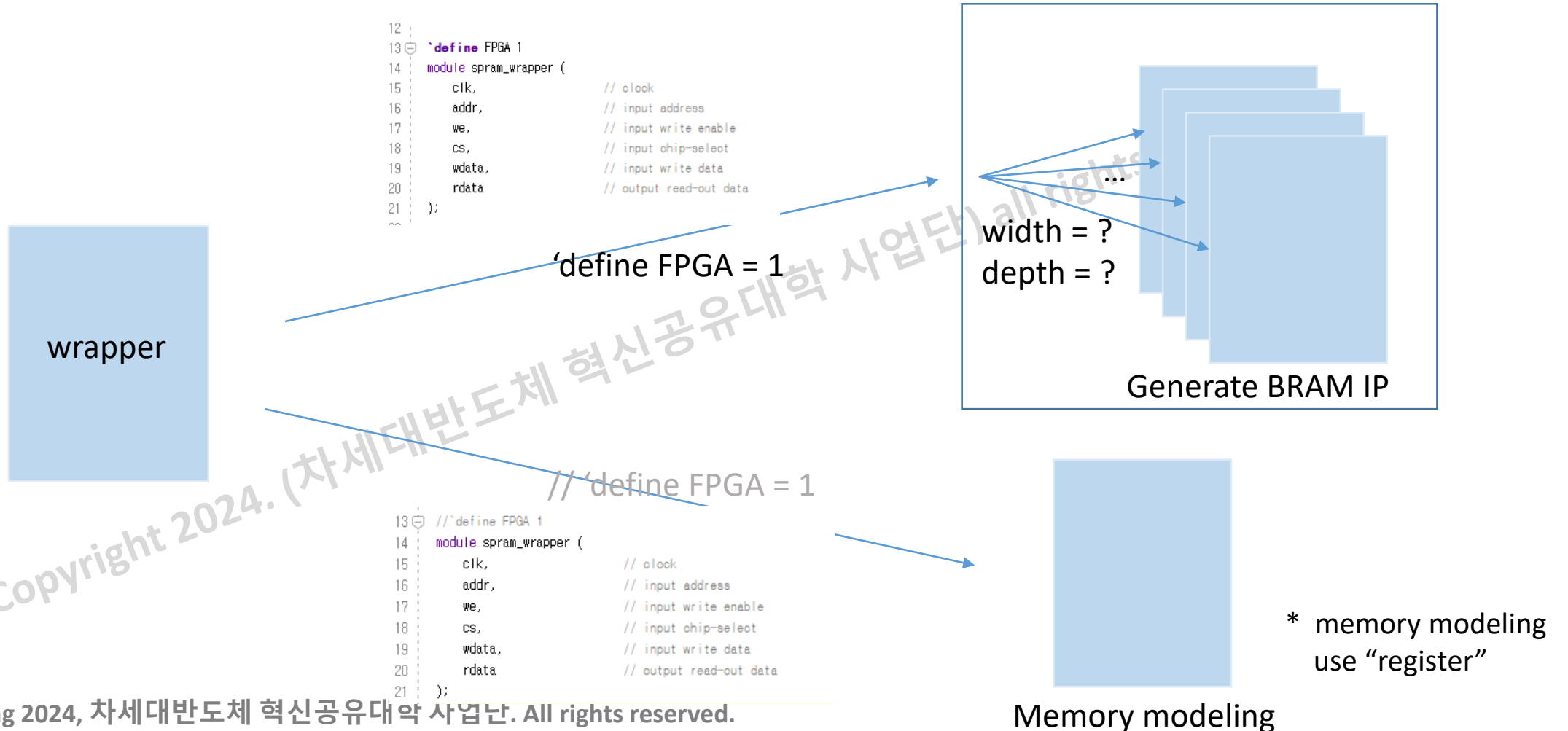
The screenshot shows the 'Component Name' field set to 'spram_6240x32'. Below this are tabs for 'Basic', 'Port A Options', 'Other Options', and 'Summary'. The 'Basic' tab is selected. Under 'Pipeline Stages within Mux', the value is '0' and 'Mux Size: 3x1'. The 'Memory Initialization' section has three red boxes highlighting the 'Load Init File' checkbox, the 'Coe File' field (containing 'sers/main/Desktop/memory_tutorial_dpram/CONV20_W.coe') with 'Browse' and 'Edit' buttons, and the 'Fill Remaining Memory Locations' checkbox. Below this is a 'Remaining Memory Locations (Hex)' field with the value '0'. The 'Structural/UniSim Simulation Model Options' section has a 'Collision Warnings' dropdown set to 'All'. The 'Behavioral Simulation Model Options' section has two unchecked checkboxes: 'Disable Collision Warnings' and 'Disable Out of Range Warnings'.

- 9. click load init file
- 10. select the coe file
- 11. check the fill remaining memory locations
- 12. click ok button

* Memory initialization : initialize the bram data to selected "~.COE" file

Wrapper (spram_wrapper.v)

- Wrapper : wrapper chooses whether to make the BRAM IP or use the memory modeling (register)



Wrapper (spram_wrapper.v) - parameter

- Parameter
 - DW : data bit width per word
 - AW : address bit width
 - DEPTH : word length

```
//-----+
// Declare parameters
//-----+
// output parameters
parameter DW = 64;      // data bit-width per word
parameter AW = 8;       // address bit-width
parameter DEPTH = 256;  // depth, word length
parameter N_DELAY = 1;
```

Ex)

DW = 4

DEPTH = 8

$AW = \log_2 DEPTH$
= 3

0	1	0	1
0	0	0	1
0	1	0	1
1	1	0	1
0	0	0	1
0	0	1	1
0	1	0	1
0	0	0	1

Wrapper (spram_wrapper.v) - parameter

- Modify Parameter : you can customize the BRAM by modifying parameters
- 1. using BRAM IP (define FPGA = 1)
 - Set parameter
 - Add the proper "else if ~" code
 - Make BRAM IP
- 2. using memory modeling
 - Just set parameter

```
// output parameters  
parameter DW = 32;  
parameter AW = 13;  
parameter DEPTH = 6240;  
parameter N_DELAY = 1;
```

Wrapper (spram_wrapper.v) - parameter

- Modify Parameter : you can customize the BRAM by modifying parameters

- 1. using BRAM IP (define FPGA = 1)

- Set parameter
- Add the proper "else if ~" code
- Make BRAM IP

- 2. using memory modeling

- Just set parameter

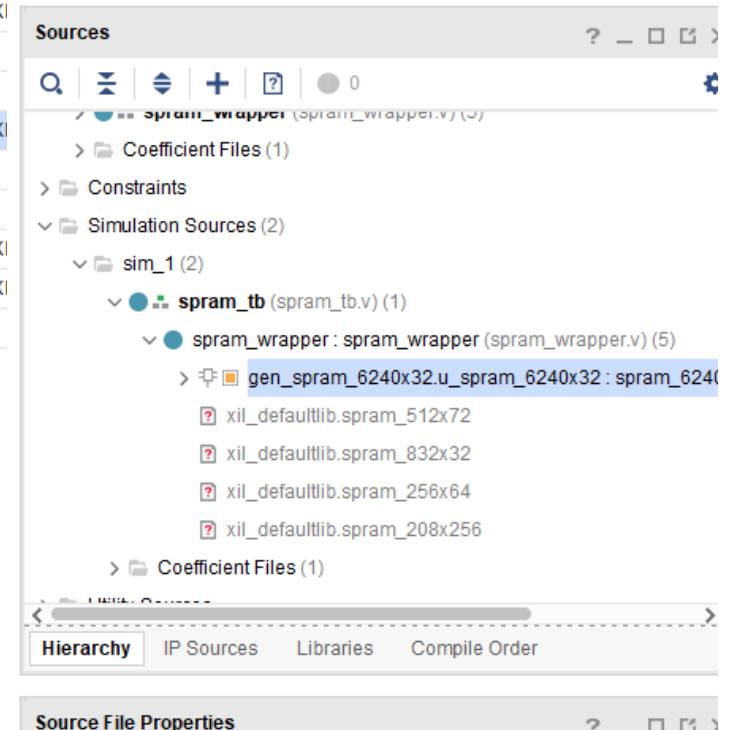
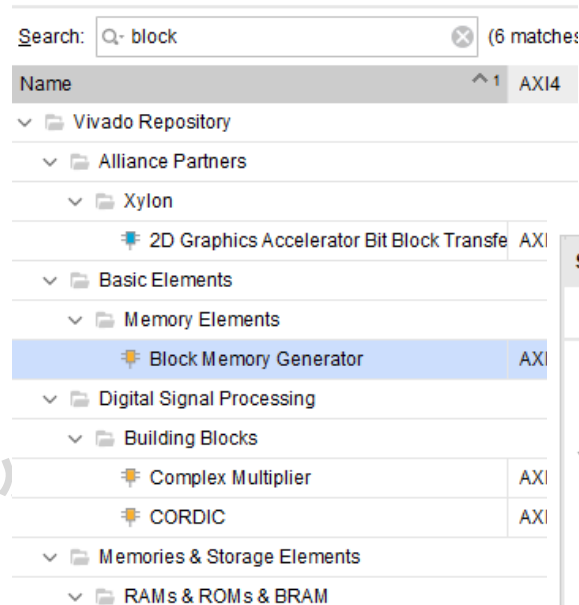
```
else if((DEPTH == 208) && (DW == 256)) begin: gen_spram_208x256
    spram_208x256 u_spram_208x256(
        // write
        .clka(clk),
        .ena(cs),
        .wea(we),
        .addra(addr),
        .dina(wdata),
        // read-out
        .douta(rdata)
    );
end
```

```
else if((DEPTH == 6240) && (DW == 32)) begin: gen_spram_6240x32
    spram_6240x32 u_spram_6240x32(
        // write
        .clka(clk),
        .ena(cs),
        .wea(we),
        .addra(addr),
        .dina(wdata),
        // read-out
        .douta(rdata)
    );
end
```

Wrapper (spram_wrapper.v) - parameter

- Modify Parameter : you can customize the BRAM by modifying parameters

- 1. using BRAM IP (define FPGA = 1)
 - Set parameter
 - Add the proper "else if ~" code
 - Make BRAM IP
- 2. using memory modeling
 - Just set parameter



Wrapper (spram_wrapper.v) - parameter

- Modify Parameter : you can customize the BRAM by modifying parameters
- 1. using BRAM IP (define FPGA = 1)
 - Set parameter
 - Add the proper "else if ~" code
 - Make BRAM IP
- 2. using memory modeling
 - Just set parameter

```
// output parameters  
parameter DW = 32;  
parameter AW = 13;  
parameter DEPTH = 6240;  
parameter N_DELAY = 1;
```


Explanation - spram_wrapper.v

- Memory modeling : just operate like BRAM IP but it's actually register

```
//-----
// Memory modeling
//
```

```
reg [DW-1 : 0] mem[0:DEPTH-1]; // Memory cell
```

Make “mem” register (data storage)

```
// Write
```

```
always @(posedge clk) begin
    if(cs && we) mem[addr] <= wdata;
end
```

Data write to “mem”

```
// Read
```

```
generate
    if(N_DELAY == 1) begin: gen_delay_1
        always @(posedge clk)
            if (cs && !(|we)) rdata_o <= mem[addr];

        assign rdata = rdata_o;
    end
```

Data read from “mem”
: 1 cycle delay

```
else begin: gen_delay_n
    reg [N_DELAY+DW-1:0] rdata_r;

    always @(posedge clk)
        if (cs && !(|we)) rdata_r[0+DW+:DW] <= mem[addr];
```

Data read from “mem”
: more than 1 cycle delay
(not used for this example! Just ignore)

```
always @(posedge clk) begin: delay
    integer i;
    for(i = 0; i < N_DELAY-1; i = i+1)
        if(cs && !(|we))
            rdata_r[(i+1)*DW+:DW] <= rdata_r[i*DW+:DW];
end
```

Support N cycle data delay
: 1 cycle delay for this example

```
assign rdata = rdata_r[(N_DELAY-1)*DW+:DW];
```

Assign “mem” output port

```
end
```

```
endgenerate
```

Test bench (spram_tb.v)

- Test bench for bram operation
 - Operation
 - 1. Read the initialized value (initialized by COE file)
 - 2. Write the test data set (width = 32, depth = 16)
 - 3. Read the changed value

```
// test data set: width = DW , depth = 16
```

```
test_data[0] = 32'h00000000;  
test_data[1] = 32'h11111111;  
test_data[2] = 32'h22222222;  
test_data[3] = 32'h33333333;  
test_data[4] = 32'h44444444;  
test_data[5] = 32'h55555555;  
test_data[6] = 32'h66666666;  
test_data[7] = 32'h77777777;  
test_data[8] = 32'h88888888;  
test_data[9] = 32'h99999999;  
test_data[10] = 32'haaaaaaaa;  
test_data[11] = 32'hbbbbbbbb;  
test_data[12] = 32'hcccccccc;  
test_data[13] = 32'hdddddddd;  
test_data[14] = 32'hEEEEEEEE;  
test_data[15] = 32'hffffffff;
```

write

read

```
// ----- read operation ----- //  
for(i=0;i<DEPTH;i=i+1) begin  
    #(8*CLK_HALF_CYCLE)  
    cs = 1'b1;  
    addr = i;  
    #(4*CLK_HALF_CYCLE)  
    cs = 1'b0;  
end
```

```
// ----- write operation ----- //  
for(i=0;i<16;i=i+1) begin  
    #(8*CLK_HALF_CYCLE)  
    cs = 1'b1;  
    we = 1'b1;  
    addr = i;  
    wdata = test_data[i];  
    #(2*CLK_HALF_CYCLE)  
    cs = 1'b0;  
    we = 1'b0;  
end
```

```
// ----- read operation ----- //  
for(i=0;i<DEPTH;i=i+1) begin  
    #(8*CLK_HALF_CYCLE)  
    cs = 1'b1;  
    addr = i;  
    #(4*CLK_HALF_CYCLE)  
    cs = 1'b0;  
end
```

Test bench (spram_tb.v) - waveform

- 1. Read the initialized value (initialized by COE file)

```
// ----- read operation ----- //
```

```
for(i=0; i<DEPTH; i=i+1) begin
```

```
    #(8*CLK_HALF_CYCLE)
```

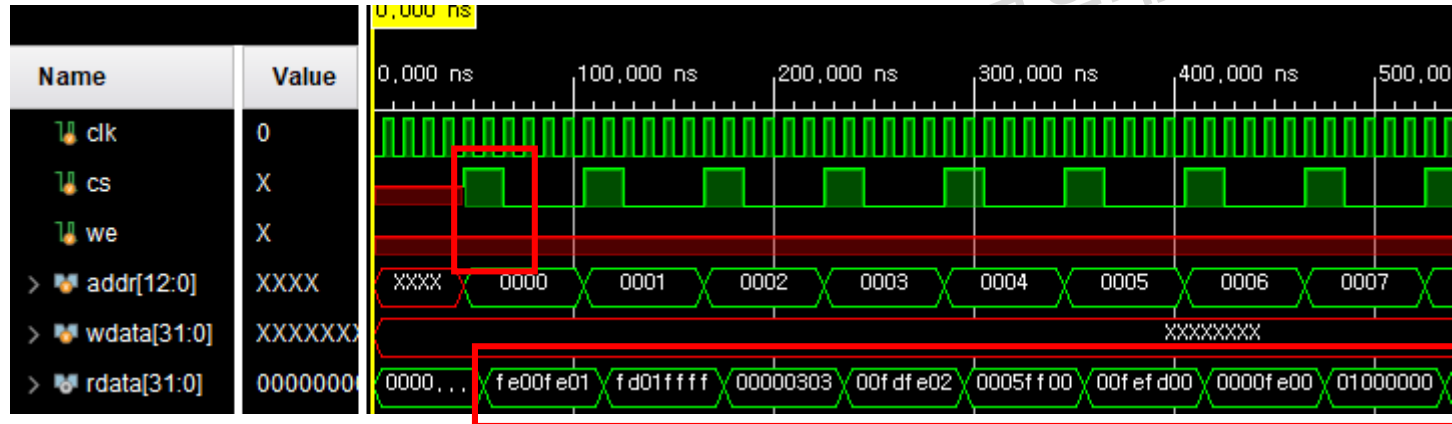
```
    cs = 1'b1;
```

```
    addr = i;
```

```
    #(4*CLK_HALF_CYCLE)
```

```
    cs = 1'b0;
```

```
end
```



CONV20_W.coe

```
1 memory_initialization_radix=16;
```

```
2 memory_initialization_vector=
```

```
3 fe00fe01 ,
```

```
4 fd01ffff ,
```

```
5 00000303 ,
```

```
6 00fdfe02 ,
```

```
7 0005ff00 ,
```

```
8 00fe0000 ,
```

```
9 0000fe00 ,
```

```
10 01000000 ,
```

```
11 fdff00fb ,
```

```
12 01feffff ,
```

```
13 ffff0201 ,
```

```
14 00fa0001 ,
```

```
15 010400fd ,
```

```
16 0500ff00 ,
```

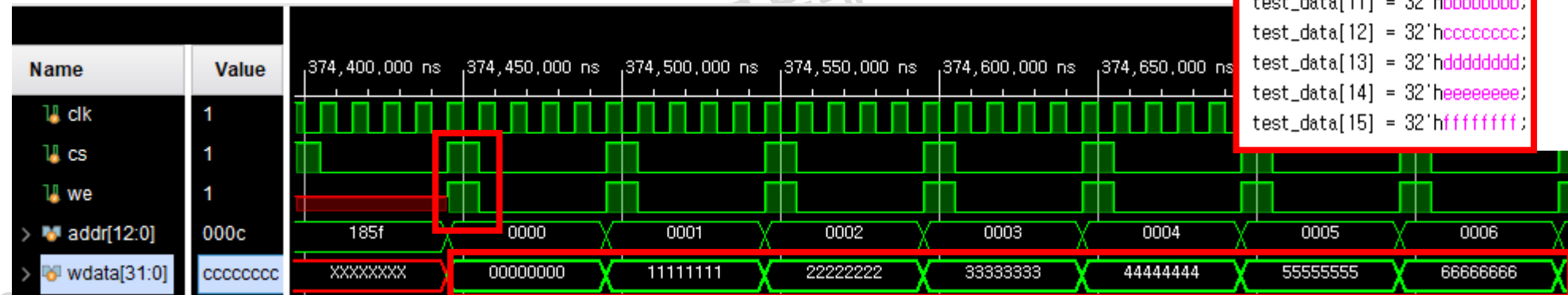
Test bench (spram_tb.v) - waveform

- 2. Write the test data set (width = 32, depth = 16);

```
// ----- write operation ----- //  
for(i=0;i<16;i=i+1) begin  
    #(8*CLK_HALF_CYCLE)  
    cs = 1'b1;  
    we = 1'b1;  
    addr = i;  
    wdata = test_data[i];  
    #(2*CLK_HALF_CYCLE)  
    cs = 1'b0;  
    we = 1'b0;  
end
```

```
// test data set: width = 32, depth = 16
```

```
test_data[0] = 32'h00000000;  
test_data[1] = 32'h11111111;  
test_data[2] = 32'h22222222;  
test_data[3] = 32'h33333333;  
test_data[4] = 32'h44444444;  
test_data[5] = 32'h55555555;  
test_data[6] = 32'h66666666;  
test_data[7] = 32'h77777777;  
test_data[8] = 32'h88888888;  
test_data[9] = 32'h99999999;  
test_data[10] = 32'haaaaaaaa;  
test_data[11] = 32'hbbbbbbbb;  
test_data[12] = 32'hcccccccc;  
test_data[13] = 32'hdddddddd;  
test_data[14] = 32'hEEEEEEEE;  
test_data[15] = 32'hffffffff;
```



Test bench (spram_tb.v) - waveform

- 3. Read the changed value

```
// ----- read operation ----- //
```

```
for(i=0;i<DEPTH;i=i+1) begin
```

```
    #(8*CLK_HALF_CYCLE)
```

```
    cs = 1'b1;
```

```
    addr = i;
```

```
    #(4*CLK_HALF_CYCLE)
```

```
    cs = 1'b0;
```

```
end
```

```
// test_data set: width = DW, depth = 16
```

```
test_data[0] = 32'h00000000;
```

```
test_data[1] = 32'h11111111;
```

```
test_data[2] = 32'h22222222;
```

```
test_data[3] = 32'h33333333;
```

```
test_data[4] = 32'h44444444;
```

```
test_data[5] = 32'h55555555;
```

```
test_data[6] = 32'h66666666;
```

```
test_data[7] = 32'h77777777;
```

```
test_data[8] = 32'h88888888;
```

```
test_data[9] = 32'h99999999;
```

```
test_data[10] = 32'haaaaaaaa;
```

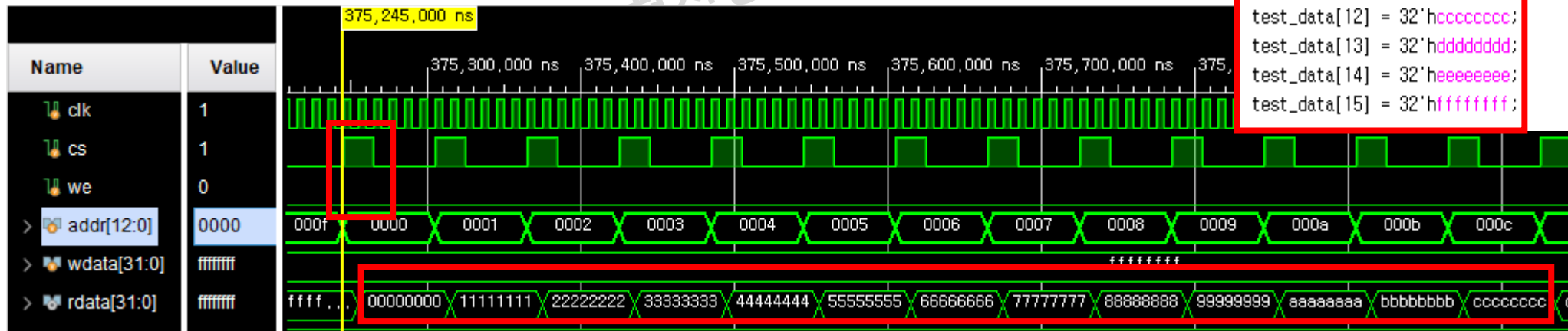
```
test_data[11] = 32'hbbbbbbbb;
```

```
test_data[12] = 32'hcccccccc;
```

```
test_data[13] = 32'hdddddddd;
```

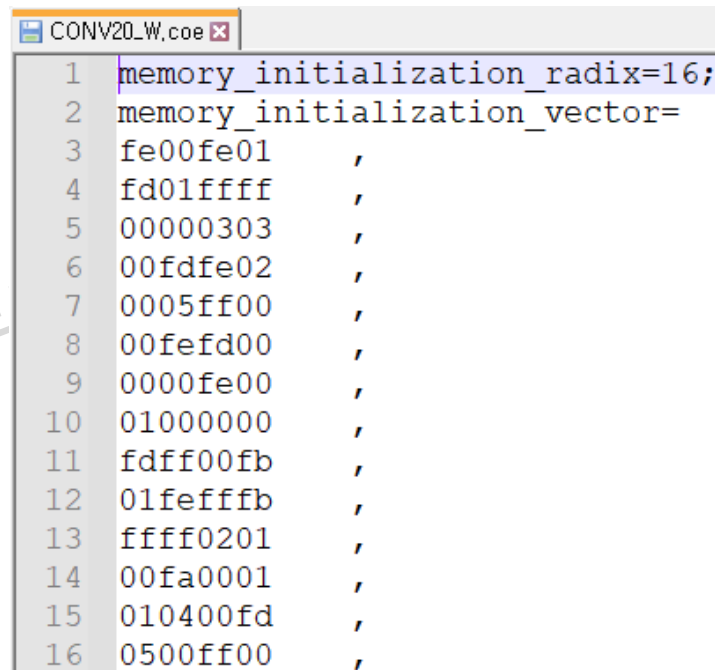
```
test_data[14] = 32'hEEEEEEEE;
```

```
test_data[15] = 32'hffffffff;
```



Lab 2: dual-port RAM

- Lab 2: Dual-port RAM (spram): 16 bits per word, 6240 words
 - Use a dual-port RAM
 - Generate a dedicated RAM using Vivado IP generator
 - Initialize a RAM using a coe file
 - Test bench



```
1 memory_initialization_radix=16;  
2 memory_initialization_vector=  
3 fe00fe01 ,  
4 fd01ffff ,  
5 00000303 ,  
6 00fdfe02 ,  
7 0005ff00 ,  
8 00fe0d00 ,  
9 0000fe00 ,  
10 01000000 ,  
11 fdff00fb ,  
12 01feffff ,  
13 ffff0201 ,  
14 00fa0001 ,  
15 010400fd ,  
16 0500ff00 ,
```


Double-port block RAM (dpram)

- Ports

- Port A : Write only

- Clock (clka)
 - Read/Write enable (ena)
 - Write enable (wea)
 - Address (addra)
 - Write data (dina)

- Port B : Read only

- Clock (clkb)
 - Read/write enable (enb)
 - Address (addrb)
 - Read data (doutb)

- Read operation

- If(enb)
 - $doutb \leq mem[addrb]$

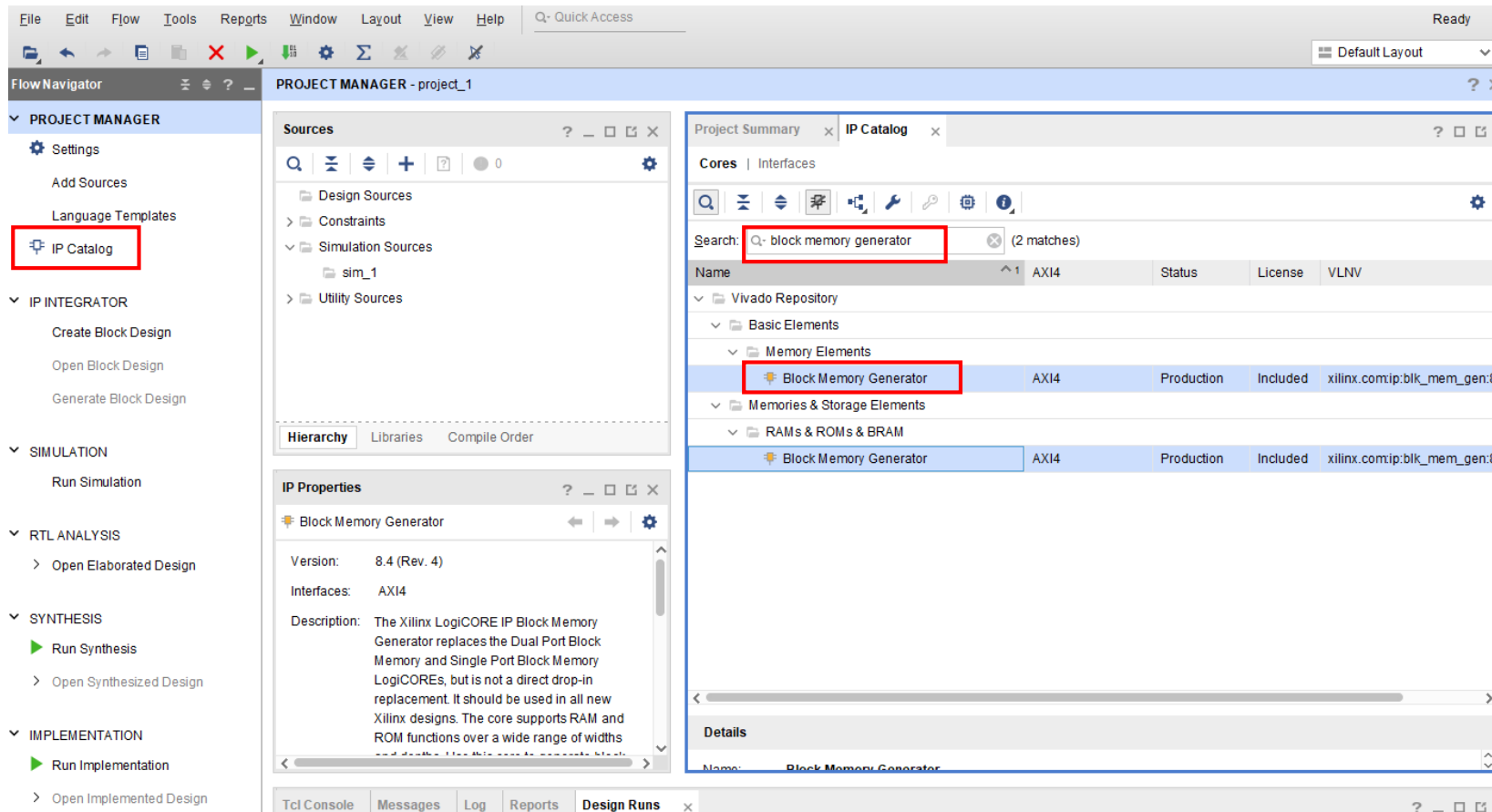
- Write operation

- If(ena & wea)
 - $mem[addra] \leq dina$

- * Double-port

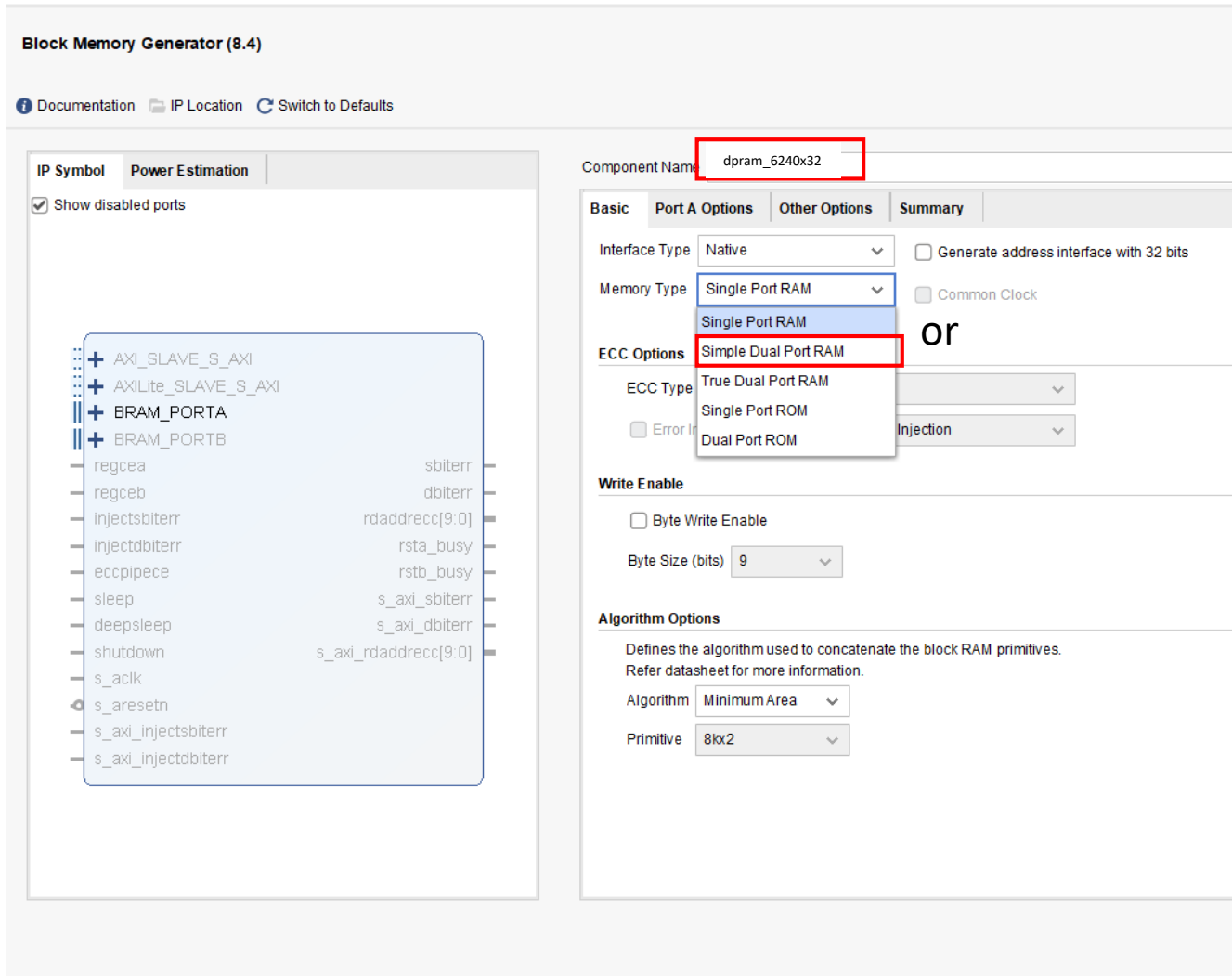
- Allow read/write operations at the same time

FPGA: How to make BRAM IP?



- 1. click IP catalog
- 2. search the "block memory generator"
- 3. double click the "block memory generator"

How to make BRAM IP?



- 4. Enter the same name that you declare on your code

- 5. Select the memory type
For reference, in given file, you will use Single Port Ram (spram) and Dual Port Ram(dpram)

Double-port block RAM (dpram_tb.v)

- About wrapper, test bench... : It's almost same with spram's
- Point
 - dpram allow read/write operations at the same time!
 - dpram write the data through port A!
 - dpram read the data through port B!

These read & write can operate at the same time, but what happen?

```
// ----- read operation ----- //
```

```
for(i=0; i<DEPTH; i=i+1) begin
```

```
    #(8*CLK_HALF_CYCLE)
```

```
    enb = 1'b1;
```

```
    addrb = i;
```

```
    #(4*CLK_HALF_CYCLE)
```

```
    enb = 1'b0;
```

```
end
```

```
// ----- write operation ----- //
```

```
for(i=0; i<16; i=i+1) begin
```

```
    #(8*CLK_HALF_CYCLE)
```

```
    ena = 1'b1;
```

```
    wea = 1'b1;
```

```
    addra = i;
```

```
    dia = test_data[i];
```

```
    #(2*CLK_HALF_CYCLE)
```

```
    ena = 1'b0;
```

```
    wea = 1'b0;
```

```
end
```

```
// ----- read operation ----- //
```

```
for(i=0; i<DEPTH; i=i+1) begin
```

```
    #(8*CLK_HALF_CYCLE)
```

```
    enb = 1'b1;
```

```
    addrb = i;
```

```
    #(4*CLK_HALF_CYCLE)
```

```
    enb = 1'b0;
```

```
end
```