

# FIR Workbook (lab\_3)

SOC Design

# Main Content

- FIR Specifications
  - FIR function
  - AXI interface
  - SRAM timing
  - Configuration Register Access Protocol
- Testbench Specifications
  - Host Software/Testbench Programming Sequence
- Additional Requirements

You will implement

- **fir.v**
- **fir\_tb.v** (testbench – you can reference and modify from Github fir\_tb.v)

# Function specification

- Same as `course-lab_2(FIRN11stream)`
- $y[t] = \sum (h[i] * x[t - i])$

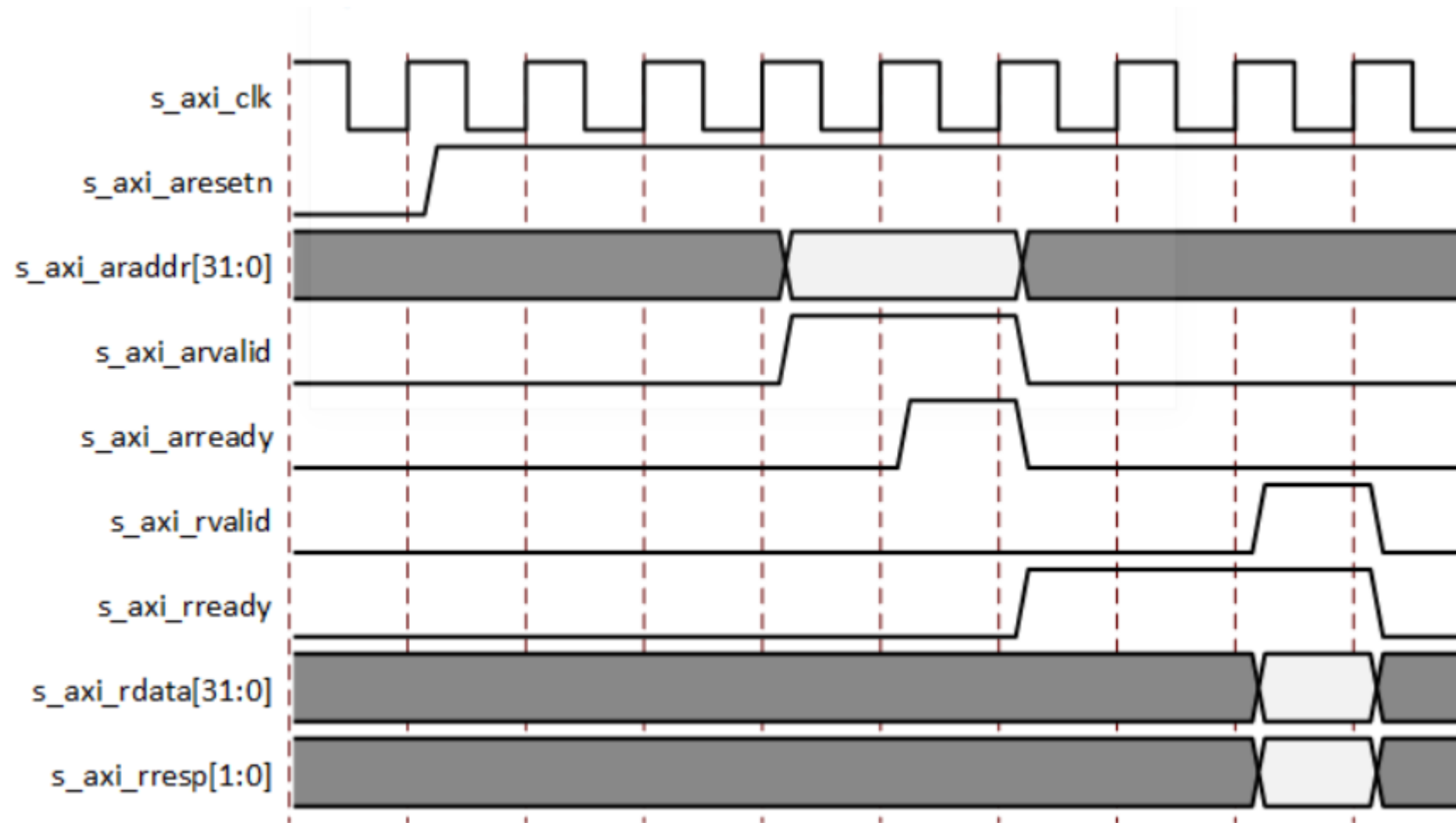
# Design specification

- Data\_Width 32
- Tape\_Num 11
- Data\_Num TBD – Based on size of data file
- Interface
  - data\_in stream ( Xn )
  - data\_out: stream ( Yn)
  - coef[Tape\_Num-1:0] axilite
  - len: axilite
  - ap\_start: axilite
  - ap\_done: axilite
- Using one Multiplier and one Adder
- Shift register implemented with SRAM (Shift\_RAM, size = 10 DW) – size = 10 DW
- Tap coefficient implemented with SRAM (Tap\_RAM = 11 DW) and initialized by axilite write
- Operation
  - ap\_start to initiate FIR engine (ap\_start valid for one clock cycle)
  - Stream-in Xn. The rate is depending on the FIR processing speed. Use axi-stream valid/ready for flow control
  - Stream out Yn, the output rate depends on FIR processing speed.

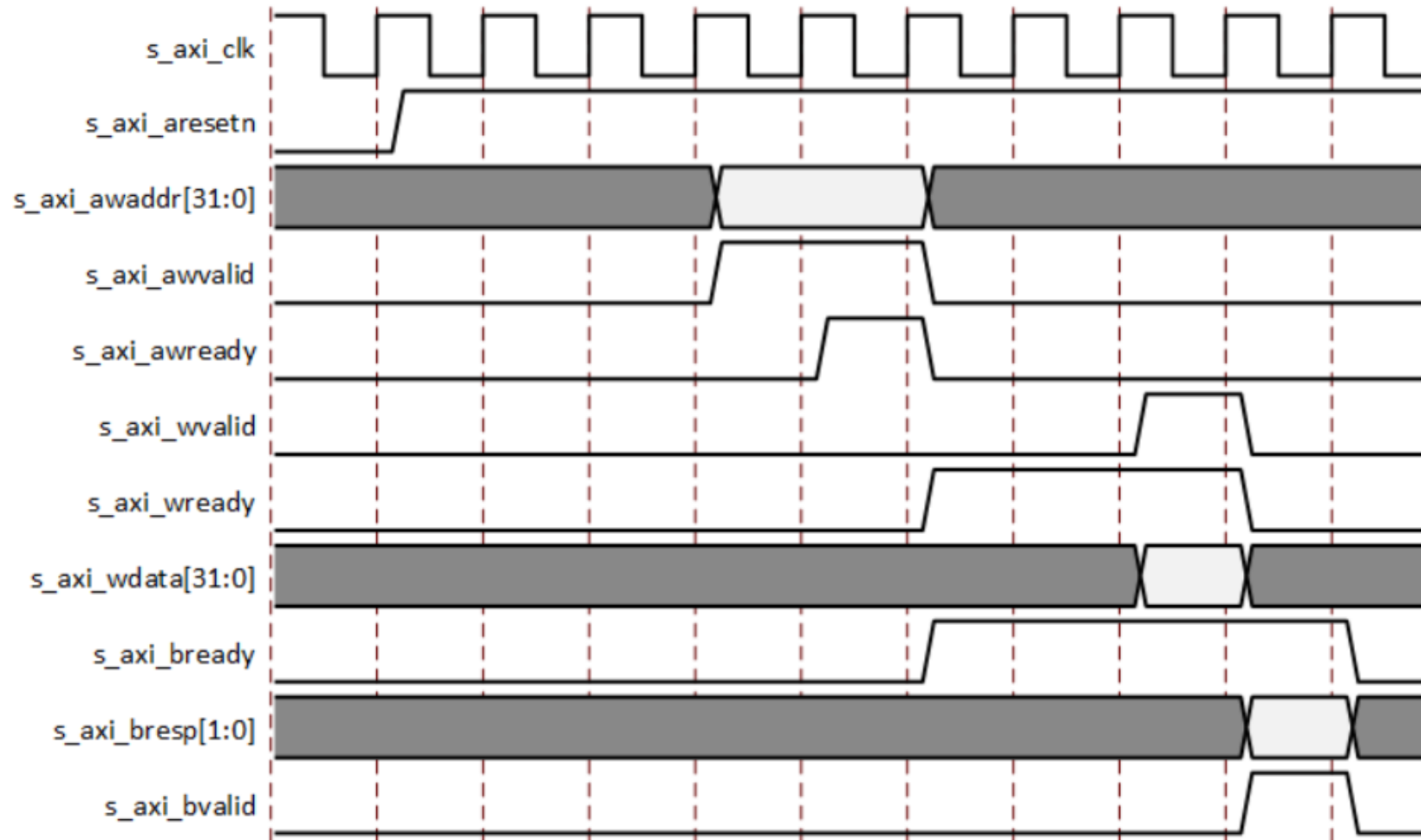
# FIR module interface (AXI-Lite, AXI-Stream)

- AXI-lite:
  - <https://www.realdigital.org/doc/a9fee931f7a172423e1ba73f66ca4081>
  - <https://docs.xilinx.com/r/en-US/pg202-mipi-dphy/AXI4-Lite-Interface>
- AXI-stream:
  - <https://developer.arm.com/documentation/ih0051/latest/>
  - <https://docs.xilinx.com/r/en-US/pg256-sdfec-integrated-block/AXI4-Stream-Interface>
- BRAM Interface: Synchronous read/write

# AXI4-Lite Read Transaction

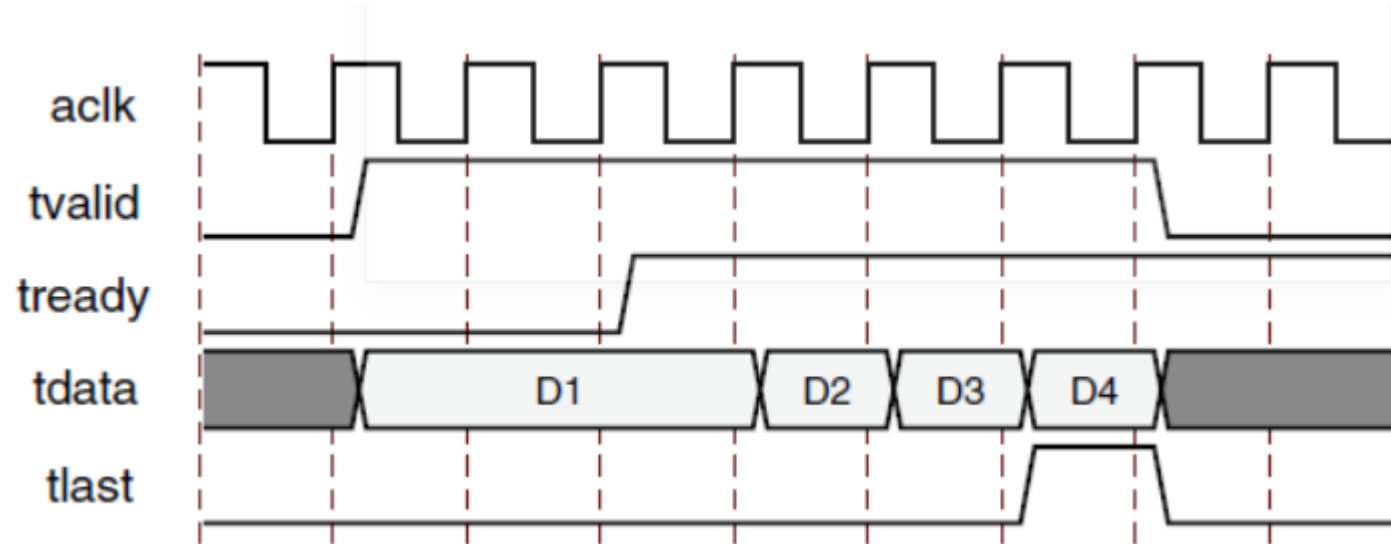
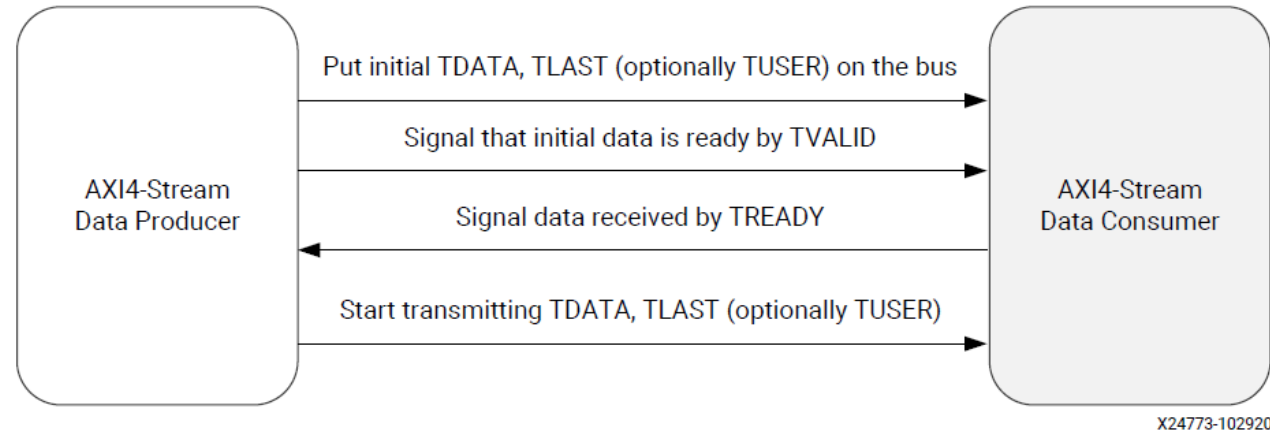


# AXI4-Lite Write Transaction





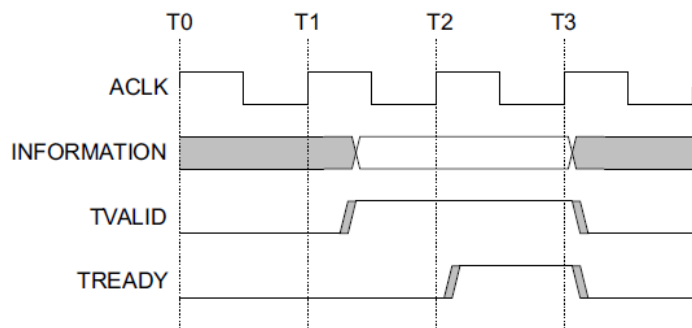
# AXI4-Stream Transfer Protocol



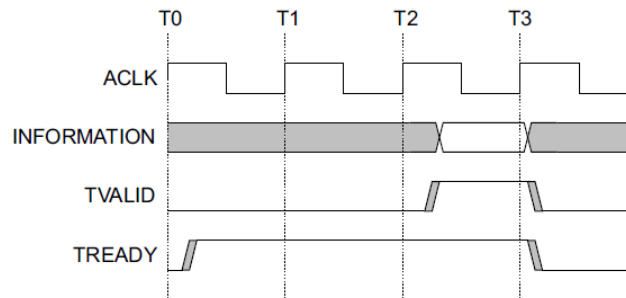
# Data Transfer Handshake : TVALID, TREADY

- For a transfer to occur, both **TVALID** and **TREADY** must be asserted
- A Transmitter is not permitted to wait until **TREADY** is asserted before asserting **TVALID**
- Once **TVALID** is asserted, it must remain asserted until the handshake occurs
- A Receiver is permitted to wait for **TVALID** to be asserted before asserting **TREADY**

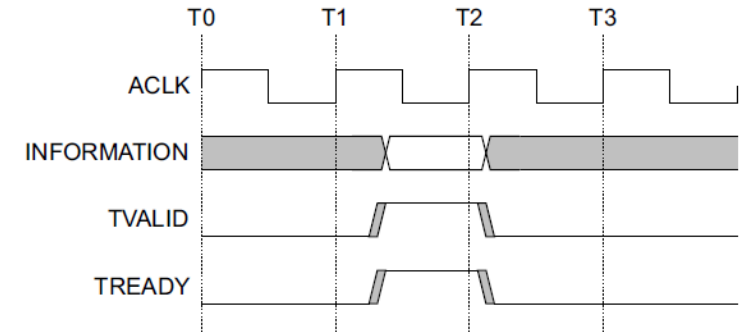
**TVALID asserted before TREADY**



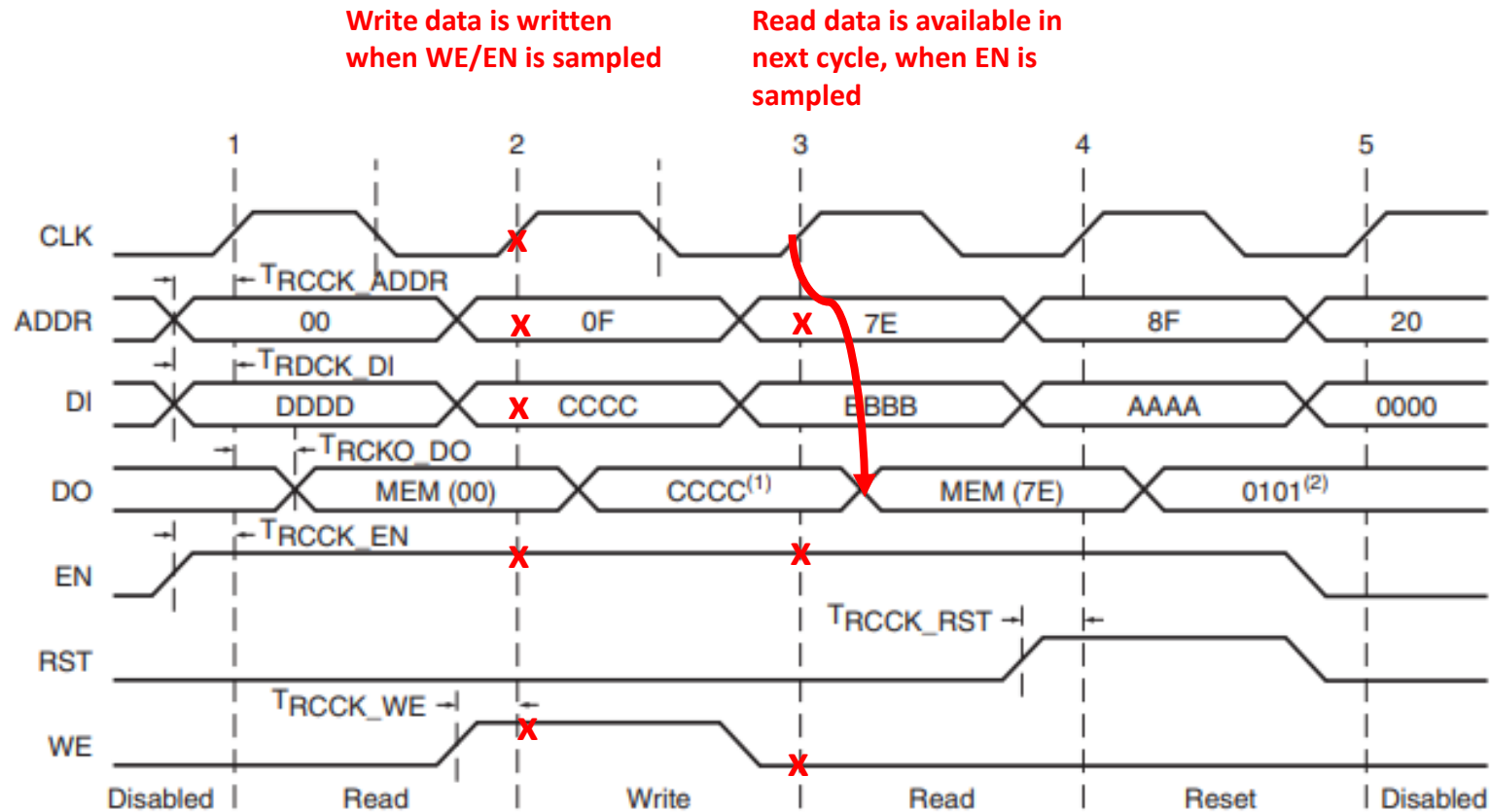
**TREADY asserted before TVALID**



**TVALID and TREADY asserted simultaneously**



# SRAM Access Timing



Note 1: Write Mode = WRITE\_FIRST

Note 2: SRVAL = 0101

UG473\_c1\_15\_052610

SRAM access has different modes, refer to <https://docs.xilinx.com/r/en-US/am007-versal-memory/Read-Operation?tocId=VRYu0HURA1U147fufYDMNQ>

# Deliver module header

- The I/O signals are listed in fir.v.
- You are requested to use simplified AXI-lite and AXI-stream protocol.

# Configuration Register Access Protocol

# Configuration Register Address map

## Address

0x00 – [0] - ap\_start (r/w)

set when ap\_start signal assert

reset, when start data transfer, i.e. 1<sup>st</sup> axi-stream data come in

[1] – ap\_done (ro) -> when FIR process all the dataset, i.e. receive tlast, and last Y is generated and transferred

[2] – ap\_idle (ro) -> indicate FIR is actively processing data

0x10-14 - data-length

0x20-FF – Tap parameters, (e.g., 0x20-24 Tap0, in sequence ...)

# ap\_start protocol and implementation

1. ap\_start is a read/write registers
2. When ap\_start is programmed one, the FIR engine starts.
3. Host Software or testbench can program ap\_start
  1. When ap\_idle is one. If ap\_start is programmed one when ap\_idle is zero, the ap\_start is not effective
  2. After data-length, tap parameters are programmed
4. ap\_start is set by software/testbench, and reset by engine
5. Engine resets ap\_start when engine is not idle, i.e. start processing data

# ap\_done protocol and implementation

1. It is read-only
2. ap\_done is reset in the following condition
  1. Reset signal is asserted
  2. When ap\_done is read, i.e. address 0 is read
3. ap\_done is asserted when engine completes last data processing and data is transferred



# ap\_idle protocol and implementation

1. ap\_idle is set to 1 when reset
2. ap\_idle is set to 0 when ap\_start is sampled
3. ap\_idle is set to 1 when FIR engine processes the last data and last data is transferred

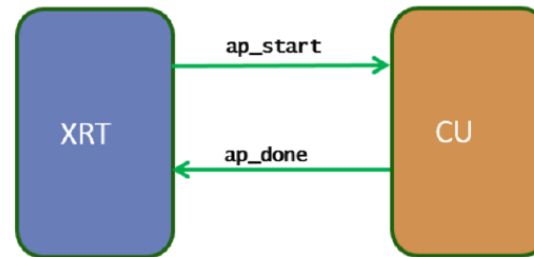
# Testbench Specification

# Testbench

- fir\_tb(Please name your top module same as **fir** for simulation)
- The testbench should keep in the same directory with Makefile
- Block-level Protocol (ap\_start, ap\_done):

AP\_CTRL\_HS (Sequential Executed Kernel)

- Host and Kernel Synchronization by
  - ap\_start
  - ap\_done
- Kernel can only be restarted (ap\_start), after it completes the current execution (ap\_done)
- Serving one execution request a time



# Host software / Testbench Programming Sequence

## Host Software / Testbench

1. Check FIR is idle, if not, wait until FIR is idle
2. Program length, and tap parameters
3. Program ap\_start -> 1
4. Fork
  1. Transmit Xn,
  2. Receive Yn
  3. Polling ap\_done
5. When ap\_done is sampled, compare Yn with golden data

## FIR Engine

Wait for ap\_start  
Set ap\_idle = 0

Process data

If reach data-length, set ap\_done

Note: Transmit Xn (stream-in), Receive Yn (stream-out) and Polling ap\_done (axilite) are running concurrently. They are using different interface and do not interfere each other

# Testbench – Develop your own testbench

## 1. Setup phase

1. Load datafile, and count # of data = data\_length
2. Program tap\_parameters and data\_length, read back and check it is correctly programmed
3. Compute Yn expected value, or load golden data into Yn buffer
4. Read and check ap\_start, ap\_idle, ap\_done are in proper state

## 2. Execution phase

1. Program ap\_start
2. Start latency timer
3. Fork the following operations, run concurrently
  1. Task: Stream\_in\_Xn
  2. Task: Stream\_out\_Yn and save into Yn buffer
  3. Task: Polling ap\_done, when ap\_done is sampled, disable tasks (stream\_in\_Xn, stream\_out\_Yn, and Polling)

## 3. Checking Phase

1. Report latency
2. Compare Yn buffer with golden data

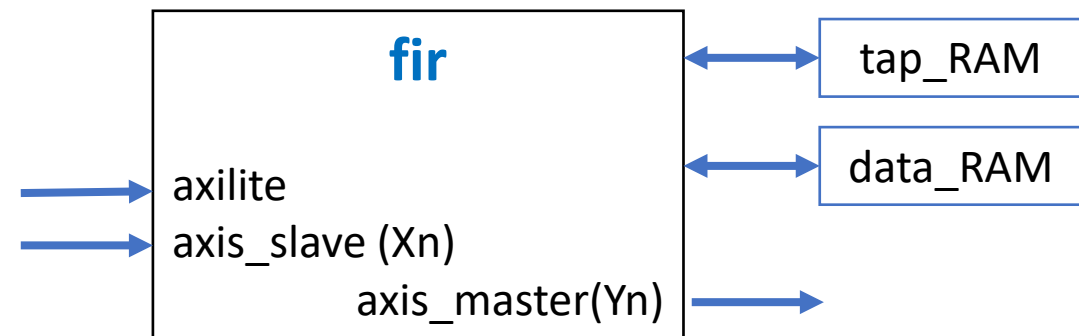
Note: You may print message to assist debugging

# Test dataset

- Samples\_triangular\_wave.dat
- out\_gold.dat

# SRAM Interface Implementation

- Refer to supplement (Use memory for ASIC flow)
- Implement SRAM without .db/.lib
- Use external SRAM (bram.v). fir.v provides ports to interface with the external SRAM. So, the fir.v can be synthesized with BRAM
- Two size of bram.v ( you choose either one to fit your design)
  - bram11.v (11 x 32) – depth 11
  - bram12.v (12 x 32) – depth 12



# Additional Requirements



# Some Design Details

- Can the FIR module fetch data continuously?
  - Can data access be pipelined with computation?
  - What if the input/output cannot process data immediately?

# Submission (1/2)

- Hierarchy:
  - StudentID\_lab3/
    - Waveform
    - Simulation.log
    - Report.pdf
    - Synthesis report
    - Github link
- Your Github link should attach the file
  - fir.v (the fir design)
  - fir\_tb.v ( the testbench )
  - Log files including : synthesis, simulation, static timing report
  - Synthesis report – area usage, Including FF, LUT (Note: there should be no BRAM because BRAM is an external model, not in the RTL design)
  - Timing Report, including slack, and max delay path
  - Waveform – show
    - Configuration write
    - ap\_start, ap\_done ( measure # of clock cycles from ap\_start to ap\_done)
    - Xn stream-in, and Yn stream-out
  - Report
- Location of design (If use vivado to design)
  - hostproject/hostproject.srscs/sources\_1/new/

# What is included in the report

- Block Diagram
  - Datapath – dataflow
  - Control signals
- Describe operation, e.g.
  - How to receive data-in and tap parameters and place into SRAM
  - How to access shiftram and tapRAM to do computation
  - How ap\_done is generated.
  - .....
- Resource usage: including FF, LUT, BRAM
- Timing Report
  - Try to synthesize the design with maximum frequency
  - Report timing on longest path, slack
- Simulation Waveform, show
  - Coefficient program, and read back
  - Data-in stream-in
  - Data-out stream-out
  - RAM access control
  - FSM

# Submission (2/2)

- Compress all above files in a single zip file named
  - StudentID\_lab3.zip (e.g., 111061545\_lab3.zip)
- Submit to HKUST Canvas
- Deadline:
  - 20% off for the late submission penalty within 3 days

Supplement

# Use Memory in ASIC Flow

# Memory Inference in ASIC

- ASIC Synthesis tool **does not infer memories from RTL** in the way FPGA synthesis tools do.
- **Use Memory Compiler**
  - Generate memory block with the specification (bitwidth, depth, # of port)
  - (.lef) – Library Exchange Format – containing placement information
  - Schematic & Netlist (for LVS and functional verification)
  - (.v) Function model (verilog) with timing check – for RTL simulation and gate-level simulation
  - (.lib/.db) – Liberty Timing File – containing Timing delay Dynamic/Static timing analysis, and synthesis
  - (.gds) – Graphical Database System – containing final layout information
- In RTL code, **explicitly instantiate the memory, and design its control signals**, e.g. Enable, read/write, address, input/output data

# Note on ASIC Implementation with SRAM

1. RTL design use memory instance directly. Need to provide sram synthesis library, either .lib, or .db (synospsy design-compiler). There is no particular inference method. (note: Xilinx FPGA can use inference)
2. If there is no sram .lib, or .db, you may put the sram outside using module ports. In this case, you can simulate with post-synthesis gate with sram behavior model. To get sram interface timing optimization, you will need to specify sram interface timing constraints, for example, output delay, input delay.



# SRAM with .db/.lib

- RTL Simulation
  - RTL code with instance of SRAM
  - Simulate with functional model
- Synthesis
  - Refer to .lib for SRAM timing/area information
  - Optimize timing for SRAM interface timing
- Post-Synthesis Gate-level Simulation
  - Post layout gate-level timing simulation
  - Use functional model with timing check (specify/endspecify)

```
// RTL module
module your_design( ...) begin

// instantiate SRAM
SRAM32X32 (CLK, WE, EN, ADDR, DI, DO)

endmodule
```

```
// Functional Model with timing check
module SRAM32X32( ... ) begin

specify // timing check
endspecify

end
```

```
// Timing model : .lib
```

# SRAM without .lib/.db

- SRAM instance could not be in RTL design for synthesis, instead, provide ports to interface with SRAM
- RTL simulation
  - Simulate with SRAM model
- Synthesis
  - Provide timing constraints (e.g. output delay, input delay ) for SRAM interface ports
- Post Synthesis with gate-level simulation
  - Simulate with SRAM functional model with timing check

```
// RTL module
module your_design(

// SRAM interface ports
    SRAM_EN,
    SRAM_ADDR,

... ) begin

// No SRAM instance
// SRAM32X32 (CLK, WE, EN, ADDR, DI, DO)

endmodule
```

```
// Functional Model with timing check
module SRAM32X32( ... ) begin

specify // timing check
endspecify

end
```

# Timing Check Tasks in Verilog

- Specify block can be used to specify setup and hold times for signals
  - **specify** and **endspecify** (Use specparam to define parameters in specify block)
- **\$setup** (data, clock edge, limit)—Displays warning message if setup timing constraint is not met
  - \$setup(d, posedge clk, 10)
- **\$hold** (clock edge, data, limit)—Displays warning message if hold timing constraint is not met
  - \$hold(posedge clk, d, 2)
- **\$width** (pulse event, limit)—Displays warning message if pulse width is shorter than limit
  - \$width(posedge clk, 20) —specify start edge of pulse
- **\$period** (pulse event, limit)—Check if period of signal is sufficiently long
  - \$period(posedge clk, 50)

# Specify Block Example

```
module d_model(  
    input    d,  
    input    clr,  
    input    clk,  
    output reg q  
);  
  
parameter Tclr = 30;  
parameter Trise = 13;  
parameter Tfall = 25;  
  
specify  
    $setup (d, posedge clk, 10); // check setup time  
    $period (posedge clk, 60);  
endspecify  
  
always @ (posedge clk, posedge clr)  
    if (clr)  
        #Tclr q <= 1'b0; // clear delay  
    else  
        if (d == 1'b1)  
            #Trise q <= 1'b1; // Tplh delay  
        else  
            #Tfall q <= 1'b0; // Tphl delay  
  
endmodule
```

```
specify  
    specparam Tsetup = 10, // minimum setup time before clk  
    Tperiod = 60; // minimum clock period  
    $setup (d, posedge clk, Tsetup); // check setup time  
    $period (posedge clk, Tperiod); // check period  
endspecify
```

SRAM Access in behavior model and Synthesizable  
Hardware Design  
ref : spiflash-vip.v v.s. spiflash.v

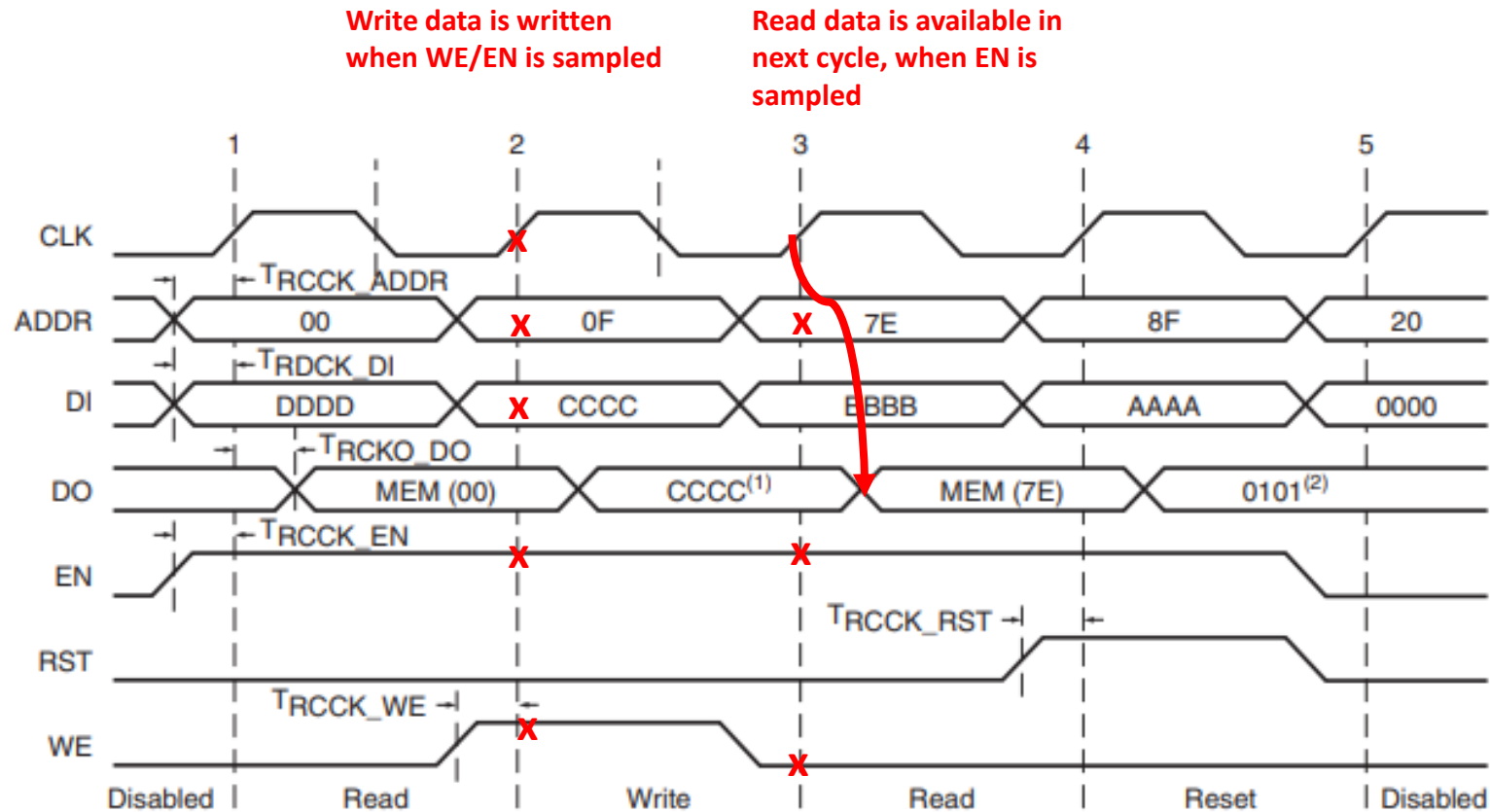
# Example: spiflash design

- spiflash-vip.v - spiflash behavior model
  - access sram as an model
- bram.v – BlockRAM behavior model
- spiflash.v
  - Adapt from spiflash-vip.v, synthesizable verilog design
  - Generate bram interface signal to access data

Reference code:

[https://github.com/bol-edu/caravel-soc\\_fpga-lab/tree/main/spiflash](https://github.com/bol-edu/caravel-soc_fpga-lab/tree/main/spiflash)

# SRAM Access Timing



UG473\_c1\_15\_052610

SRAM access has different modes, refer to <https://docs.xilinx.com/r/en-US/am007-versal-memory/Read-Operation?tocId=VRYu0HURA1U147fufYDMNQ>

# BRAM Model

**@(posedge CLK) if EN0 is sampled,  
output its memory content to Do0**

**@(posedge CLK) if WE[3:0] is  
sampled, RAM is written with Di0  
per byte**

```
module bram #( parameter FILENAME = "firmware.hex")
(
    input  wire  CLK;
    input  wire  [3:0]  WE0;
    input  wire  EN0;
    input  wire  [31:0] Di0;
    output reg  [31:0] Do0;
    input  wire  [31:0] A0
)
reg [7:0] RAM[0:4*1024*1024-1]; // Declare Memory Storage
```

```
always @(posedge CLK)
    if(EN0) begin
        Do0 <= {RAM[{A0[31:2],2'b11}],
                RAM[{A0[31:2],2'b10}],
                RAM[{A0[31:2],2'b01}],
                RAM[{A0[31:2],2'b00}]};
        if(WE0[0]) RAM[{A0[31:2],2'b00}] <= Di0[7:0];
        if(WE0[1]) RAM[{A0[31:2],2'b01}] <= Di0[15:8];
        if(WE0[2]) RAM[{A0[31:2],2'b10}] <= Di0[23:16];
        if(WE0[3]) RAM[{A0[31:2],2'b11}] <= Di0[31:24];
    end
    else
        Do0 <= 32'b0;
```

```
initial begin
    $display("Reading %s", FILENAME);
    $readmemh(FILENAME, RAM);
    $display("%s loaded into memory", FILENAME);
    $display("Memory 5 bytes = 0x%02x 0x%02x 0x%02x 0x%02x 0x%02x",
            RAM[0], RAM[1], RAM[2], RAM[3], RAM[4]);
end
```



# spiflash-vip – behavior model to access RAM

**Memory data is available the same time address is supplied. This is not feasible in the actual memory system.**

## Memory defined and initialization

```
reg [7:0] memory [0:16*1024*1024-1]; // 16MB
initial begin // memory content loaded data from file
    $readmemh(FILENAME, memory);
end
```

memory read access

```
buffer = memory[spi_addr]; // memory read
```

memory write access

```
memory[spi_addr] = buffer;
```

# spiflash.v – Synthesizable hardware design

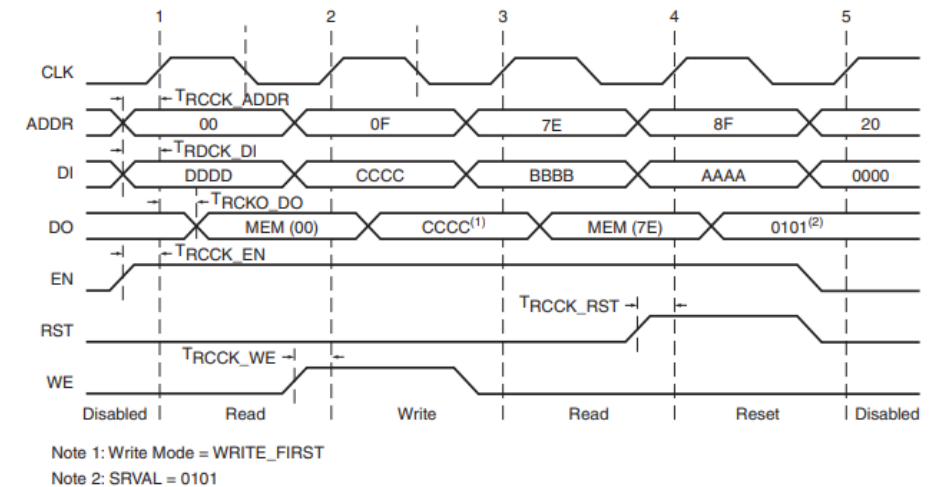
- Generate SRAM interface signals
  - Addr, EN, WEN, Din, Dout
  - **The interface signal is generated from internal control logic (FSM)** Interface signals follows the interface timing specification, e.g. **read data is available in next clock cycle**

```
// BRAM Interface
assign romcode_Addr_A = {8'b0, spi_addr};
assign romcode_Din_A = 32'b0;
assign romcode_EN_A = (bytecount >= 4);
assign romcode_WEN_A = 4'b0;
assign romcode_Clk_A = ap_clk;
assign romcode_Rst_A = ap_rst;

wire [7:0] memory;
assign memory =
    (spi_addr[1:0] == 2'b00) ? romcode_Dout_A[7:0] :
    (spi_addr[1:0] == 2'b01) ? romcode_Dout_A[15:8] :
    (spi_addr[1:0] == 2'b10) ? romcode_Dout_A[23:16] :
    romcode_Dout_A[31:24];
```

**BRAM RAM**

Addr  
EN  
WEN  
Din  
Dout  
  
Clk  
Rst



UG473\_ct\_15\_052610