Communications and Computer Engineering II

# FPGA Basis

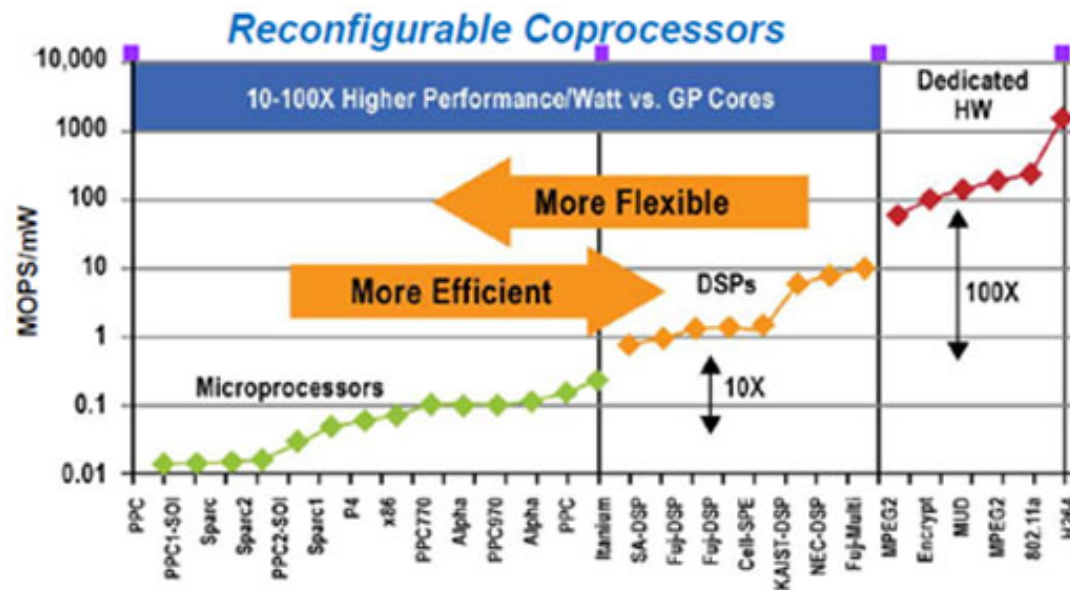Hiroki Nakahara

Tokyo Institute of Technology

# Outline

- What's FPGA?
  - Recent trends
- FPGA Architecture
- Standard FPGA Design
  - RTL (Register Transfer Level)
  - High-level synthesis (HLS)
  - System-level design
- Summary

# 1.1 FPGA

# The Dilemma: Flexibility vs. Efficiency

- FPGAs often offer the best of both worlds – replacing MPUs, DSPs, and dedicated ASSPs or ASICs

- Their on-the-fly reconfigurability helps them realize in-system logic functions that CPUs can't
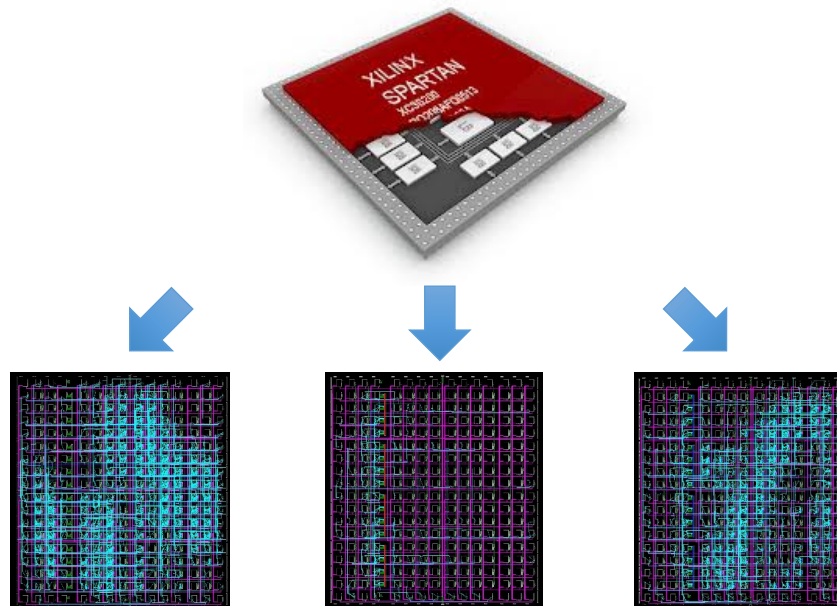


Source: "High-Performance Energy-Efficient Reconfigurable Accelerator Circuits for the Sub-45nm Era", July 2011 by Ram K. Krishnamurthy, Circuits Research Labs, Intel Corp.

# FPGA

- Reconfigurable LSI or Programmable Hardware

- Programmable Logic Array and Programmable Interconnection

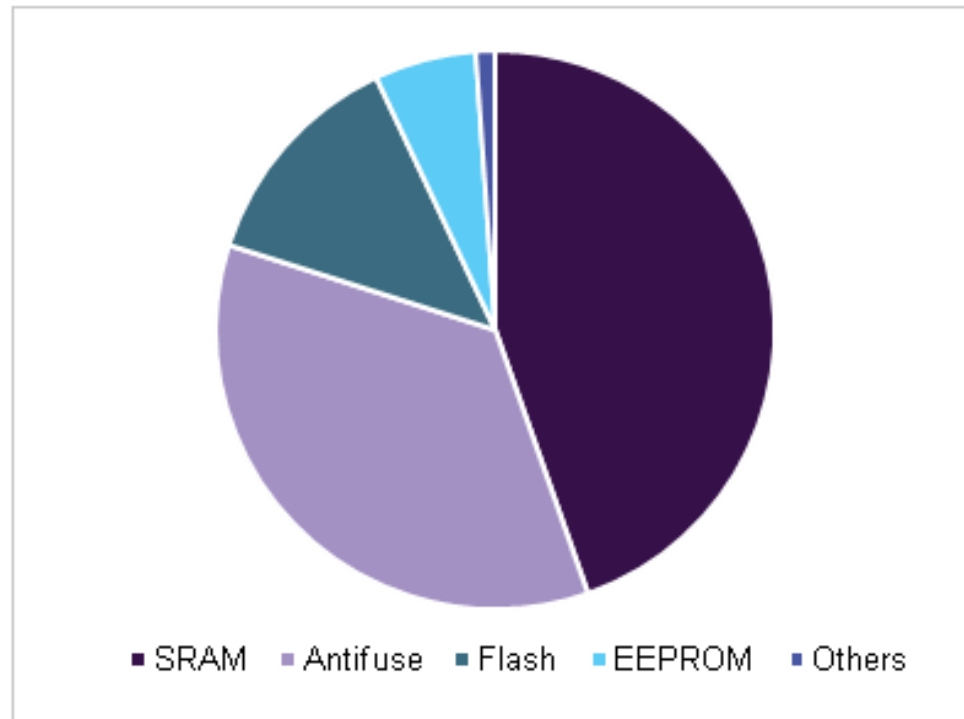- Programmed by Reconfigurable Data

# Pros. and Cons.

- Pros.
    - Short TAT(Turn-Around Time)
    - Small NRE (Non Recurrent Expense) Fee
    - Logic and Timing Design are required.
    - Full amount of IP (Intellectual Property)
- Cons.
    - Slow speed and Large Chip Area
    - High cost for volume manufacturing

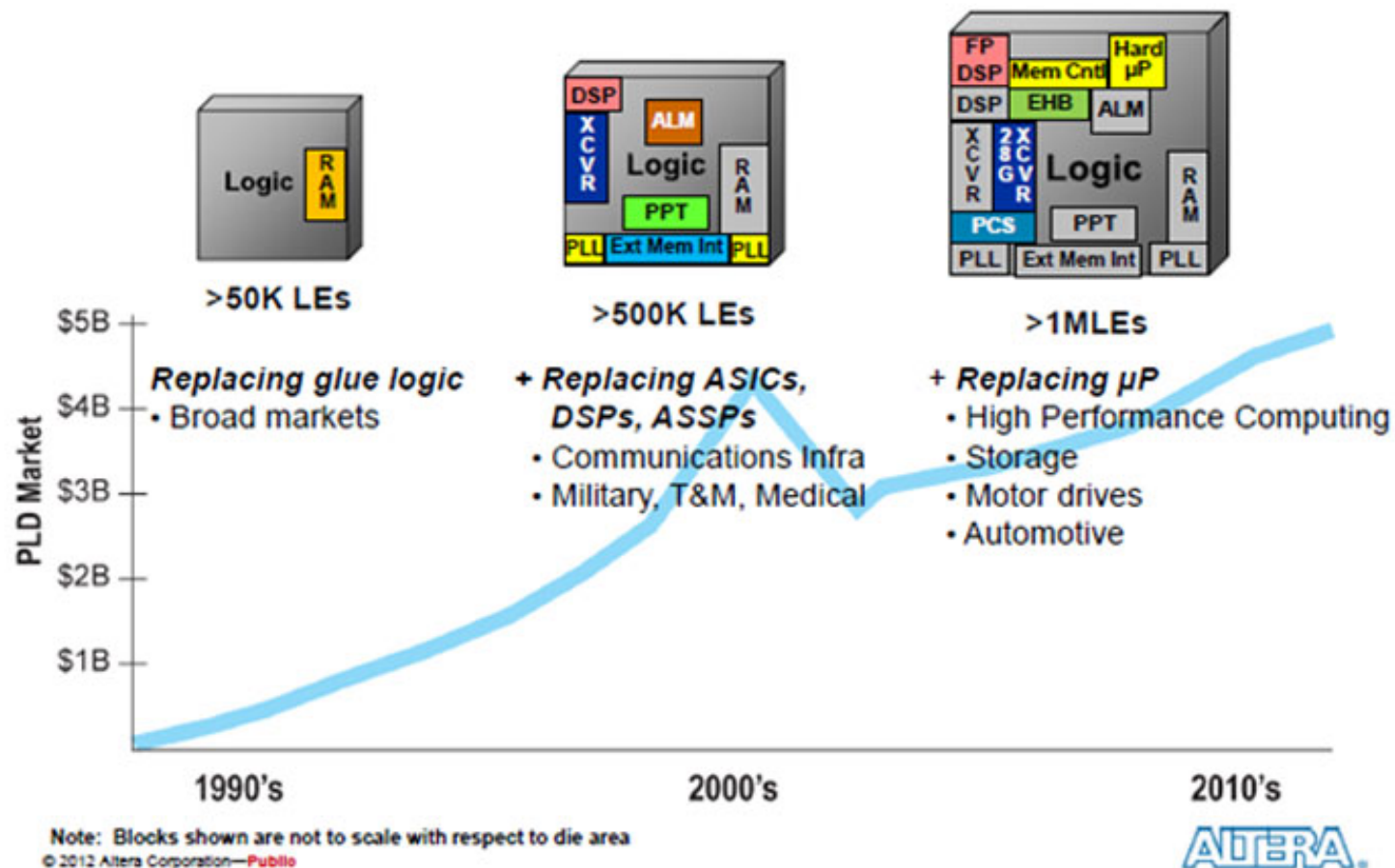# Product Type Segments

- SRAM
- Flash based
- Antifuse

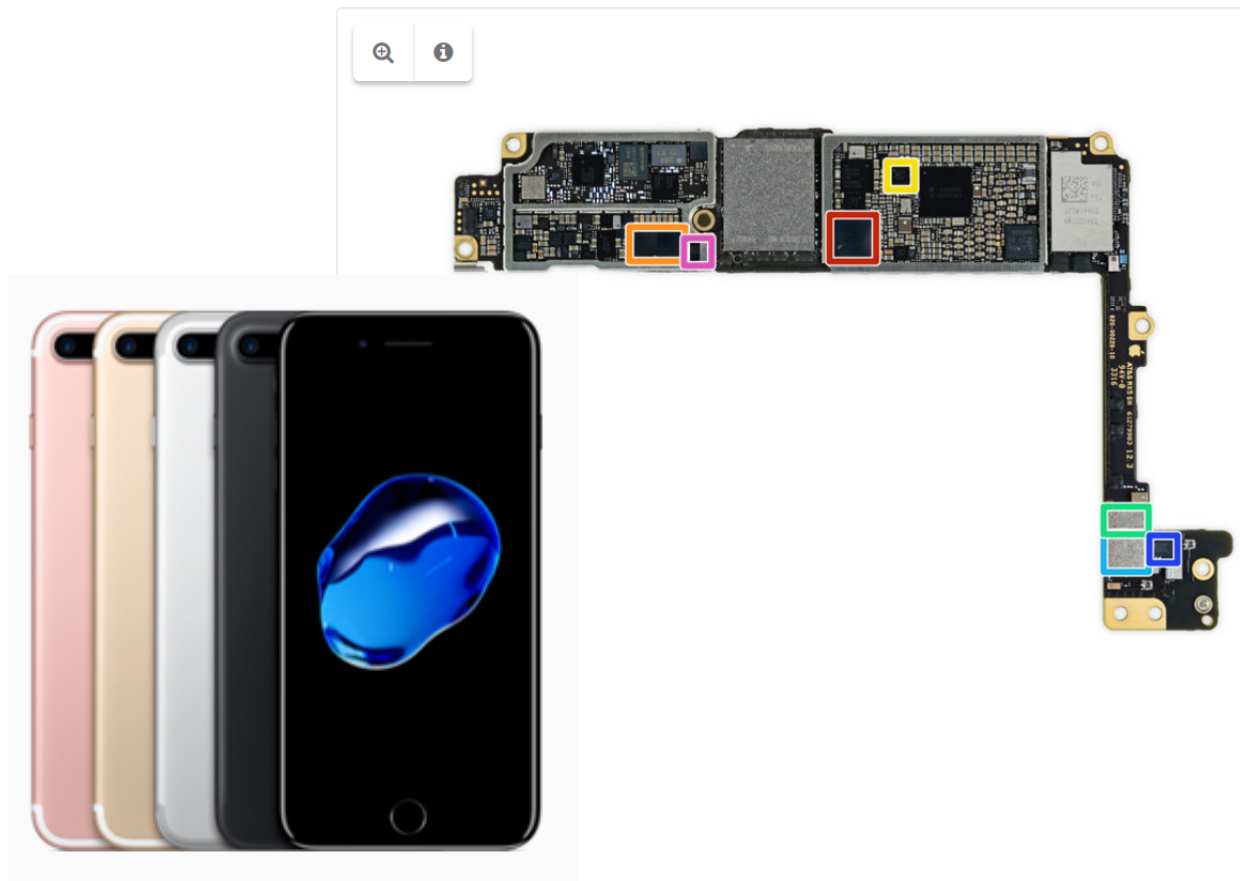**Global FPGA market share, by technology, 2015 (USD Million)**



Legend: ■ SRAM ■ Antifuse ■ Flash ■ EEPROM ■ Others

Source: https://www.grandviewresearch.com/industry-analysis/fpga-market

# FPGA Growth Trend

- 20 Years FPGAs have been swallowing up system components (by Altera, now a part of Intel)



Note: Blocks shown are not to scale with respect to die area
© 2012 Altera Corporation—Public

# Application Type Segments

- Industrial
- Automotive
- Consumer electronics
- Military & aerospace
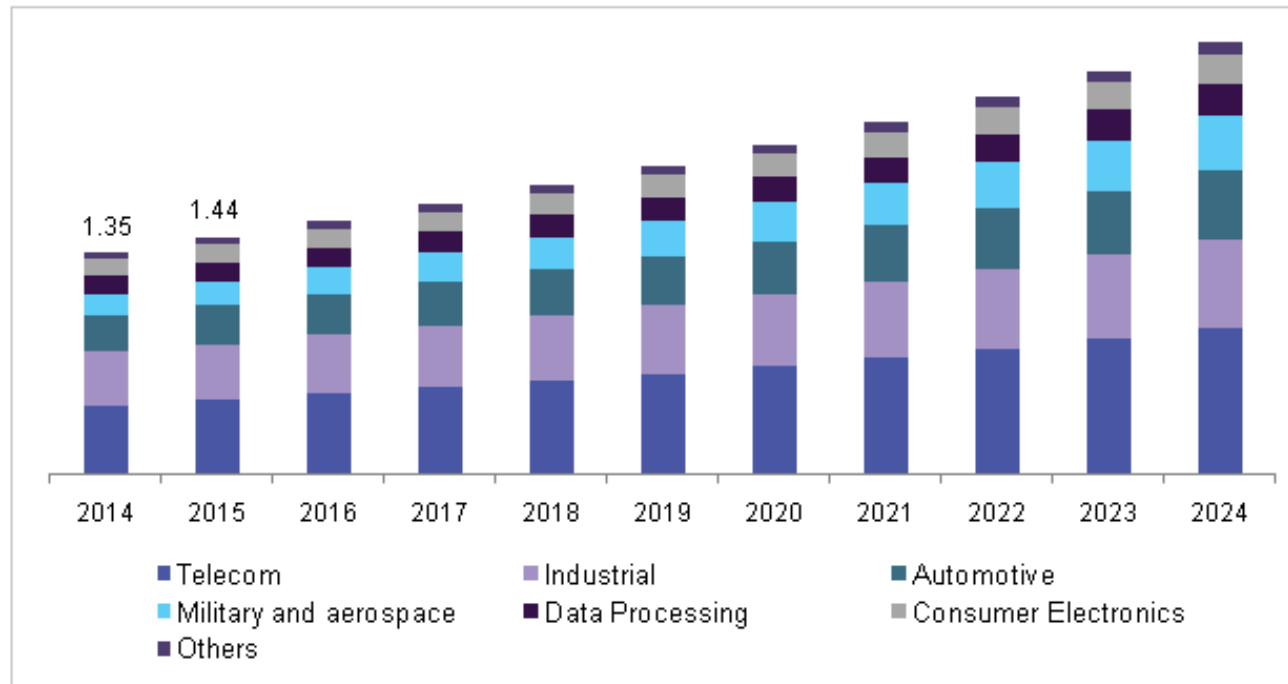- Telecom
- Data processing
- Others

# iPhone7 Plus



But wait, there are even more ICs on the back!

- Apple/Cirrus Logic 338S00105 Audio Codec
- Cirrus Logic 338S00220 Audio Amplifier(x2)
- Lattice Semiconductor ICE5LP4K
- Skyworks 13702-20 Diversity Receive Module
- Skyworks 13703-21 Diversity Receive Module
- Avago LFI630 183439
- NXP 610A38

Source: https://www.ifixit.com/Teardown/iPhone+7+Plus+Teardown/67384

# Market by Application



U.S. FPGA Market by application, 2014 - 2024 (USD Billion)

Legend: Telecom, Industrial, Automotive, Military and aerospace, Data Processing, Consumer Electronics, Others

# Market Share by Vendor

| Vendor | 2015 | | 2016 | | |
|---|---|---|---|---|---|
| | FPGA Total | Market share | FPGA Total | Market share | Growth CY15-CY16 |
| Xilinx | $2,044 | 53% | $2,167 | 53% | 6% |
| Intel (Altera) | $1,389 | 36% | $1,486 | 36% | 7% |
| Microsemi | $301 | 8% | $297 | 7% | -1% |
| Lattice | $124 | 3% | $144 | 3% | 16% |
| QuickLogic | $19 | 0% | $11 | 0% | -40% |
| Others | $2 | 0% | $2 | 0% | 0% |
| TOTAL | $3,879 | 100% | $4,112 | 100% | 6% |

Source EEtimes 3/5/2017

# FPGA vs. ASIC: Which is Best?

| anysilicon | FPGA | ASIC |
|---|---|---|
| Time to Market | Fast | Slow |
| NRE | Low | High |
| Design Flow | Simple | Complex |
| Unit Cost | High | Low |
| Performance | Medium | High |
| Power Consumption | High | Low |
| Unit Size | Medium | Low |

AnySilicon.com

Source: https://anysilicon.com/fpga-vs-asic-choose/

# Price Comparison

- ASIC→NRE: $1.5M, Unit cost: $4
- FPGA→NRE: $0, Unit cost: $8



Total Cost ASIC vs FPGA including NRE in MUSD

anysilicon.com

# Performance Comparison

- The divergence between primary programmable logic device technology and that used for ASICs has continued to grow



Source: http://archive.rtcmagazine.com/articles/view/102503

# 1.2 FPGA Architecture

# Island Style FPGA



**Configurable Logic Block (CLB)**

**I/O Block**

**Channels**

**Switch Blocks**

# Xilinx CLB

CLB

CLB

CLB

CLB

## Configurable logic block (CLB)

### Slice

Logic cell

Logic cell

### Slice

Logic cell

Logic cell

### Slice

Logic cell

Logic cell

### Slice

Logic cell

Logic cell

# Logic Cell (Xilinx Inc. XC2000)

Programmable Multiplexer

# Realization of a Logic Function

| x0 | x1 | x2 | y |
|----|----|----|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 |

| x0 | x1 | x2 | y |
|----|----|----|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 |

# LUT Structure

# Channel and Switch Block

# I/O Block

# Memory-based realizes "programmable"



| x0 | x1 | x2 | y |
|----|----|----|---|
| 0  | 0  | 0  | 0 |
| 0  | 0  | 1  | 0 |
| 0  | 1  | 0  | 0 |
| 0  | 1  | 1  | 0 |
| 1  | 0  | 0  | 0 |
| 1  | 0  | 1  | 0 |
| 1  | 1  | 0  | 0 |
| 1  | 1  | 1  | 1 |

| x0 | x1 | x2 | y |
|----|----|----|---|
| 0  | 0  | 0  | 0 |
| 0  | 0  | 1  | 1 |
| 0  | 1  | 0  | 1 |
| 0  | 1  | 1  | 1 |
| 1  | 0  | 0  | 1 |
| 1  | 0  | 1  | 1 |
| 1  | 1  | 0  | 1 |
| 1  | 1  | 1  | 1 |

どちらも
トランジスタの
スイッチ

# FPGA-world Famous Patents

US RE34363
by Ross H. Freeman
Island style architecture

US 464248
by William S. Carter
Connection architecture

# Hard macro (Dedicated Circuit)

- Standard IP cores are "Pre-built" in the FPGA



36Kb Block Memory



48bit DSP Block

# Xilinx FPGA Families

- 3D IC type
  - Virtex UltraScale+, Kintex UltraScle
- System on Chip (SoC) type
  - Zynq UltraSCale+
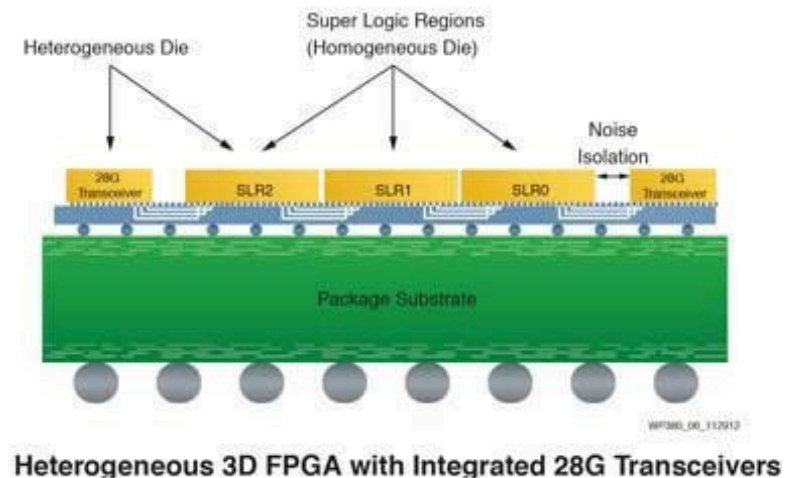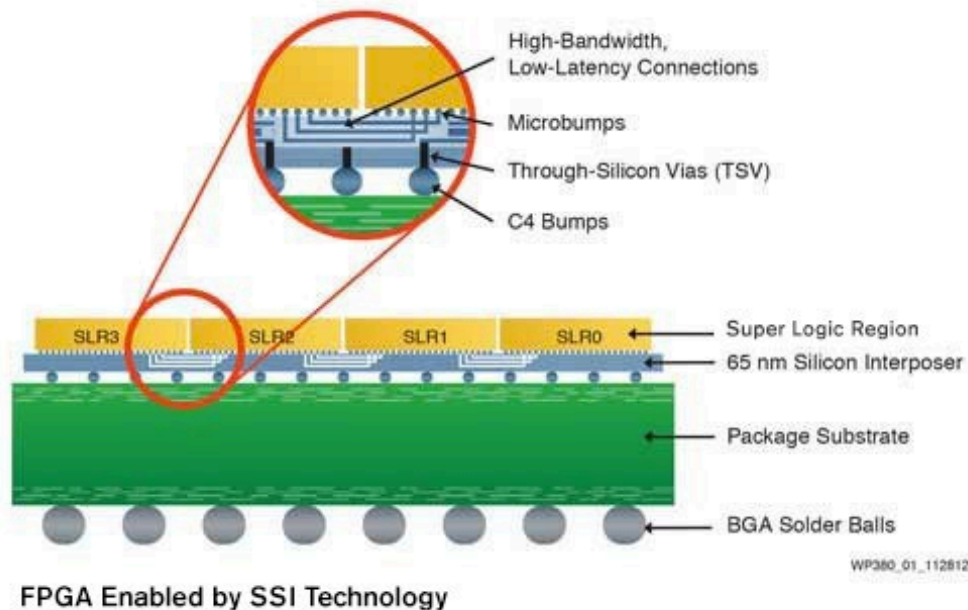- Monolithic type
  - Artix(Low-cost), Spartan (Low-cost)



First 3D FPGA: Virtex-7 2000T
Based on Stacked Silicon Interconnect

First Heterogeneous 3D FPGA: Virtex-7 H580T
Based on Stacked Silicon Interconnect

7 series | UltraScale Architecture

3D IC

SoC

FPGA

VIVADO

Programmable systems integration

FPGA die + cutting edge transceivers + wide memory + embedded SoC

Heterogeneous multi-core + FPGA processing + software environment

30-50% price/performance/watt + GT/IO Bandwidth+ IP sub-system

28nm HPL | 20nm SoC | 16nm FinFET | 10nm

Process

System-level price/performance/watt*

* System level combines unit with BOM cost

Source: http://www.xilinx.com

# Xilinx Zynq Ultra Scale+

# SSI(Stacked-Silicon Interconnect) Technology: Heterogeneous 3D FPGA

- High bandwidth connectivity across multiple dies

- Significantly increased die-to-die bandwidth per watt over multi-chip solutions

- Through-silicon via (TSV) with coarse pitch in passive silicon interposer (65nm) without transistor



FPGA Enabled by SSI Technology



Heterogeneous 3D FPGA with Integrated 28G Transceivers

Source: http://www.xilinx.com

# ACAP (Adaptive Compute Acceleration Platform)

- VERAL AI core series and prime series
  - 7nm technology
- Connected 3 elements with high-band NoC and Memory



Source: http://www.xilinx.com

# 1.3 FPGA Design

# Standard FPGA Design



**LSI Tool**

**FPGA Tool**

RTL

RTL Simulation

Logic Synthesis

Synthesis Netlist

Functional Verification

Gate Assignment LE Place and Rout

Configuration Data

HDL

# FPGA Design Flow

# How to write a HDL?

Y = X.dot(W)+B

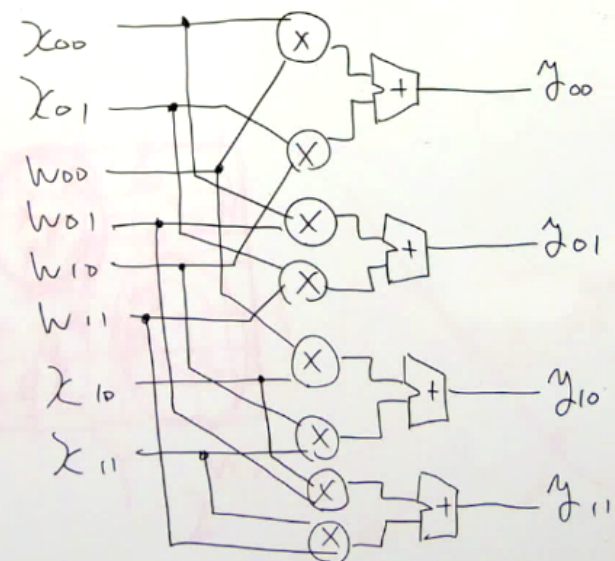

Then, you can write your HDL!

# Matrix Multiplication

$$Y = X.dot(W)$$

$\underline{Behavior}$

$$Y = \begin{pmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{pmatrix} \begin{pmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{pmatrix}$$
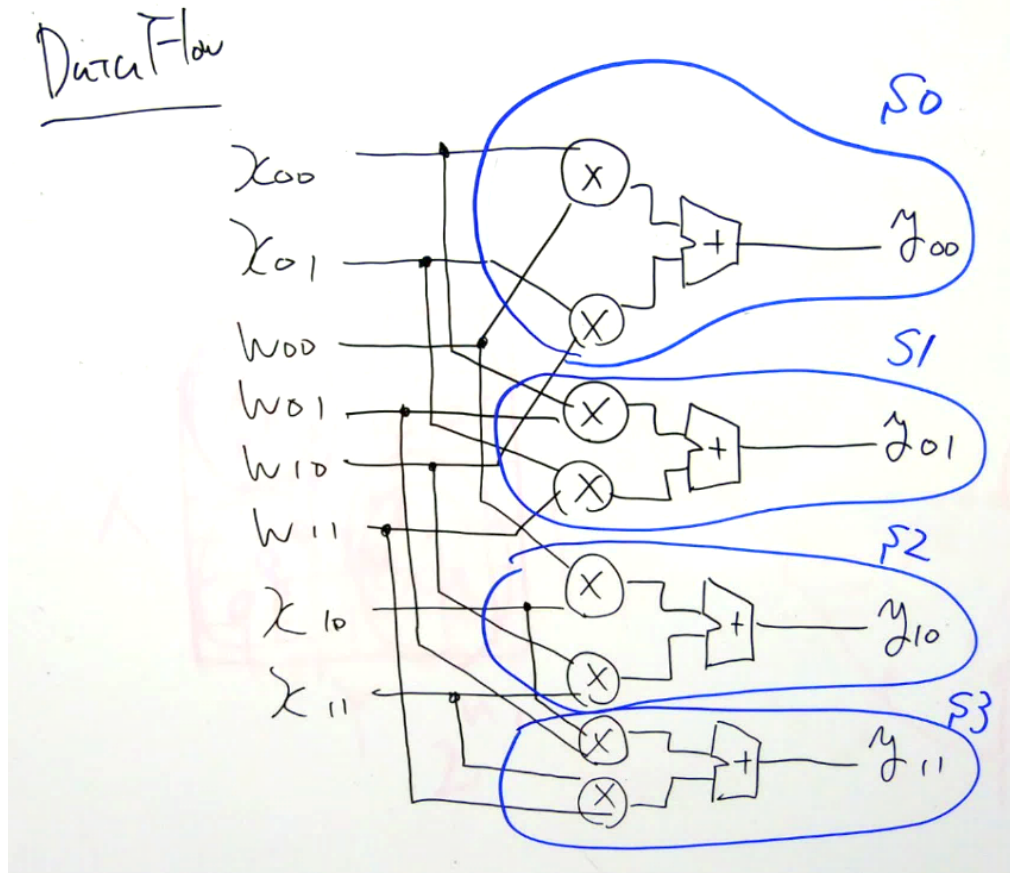
$$\begin{pmatrix} y_{00} & y_{01} \\ y_{10} & y_{11} \end{pmatrix} = \begin{pmatrix} X_{00} \cdot W_{00} + X_{01} \cdot W_{10} & X_{00} \cdot W_{01} + X_{01} \cdot W_{11} \\ X_{10} \cdot W_{00} + X_{11} \cdot W_{10} & X_{10} \cdot W_{01} + X_{11} \cdot W_{11} \end{pmatrix}$$
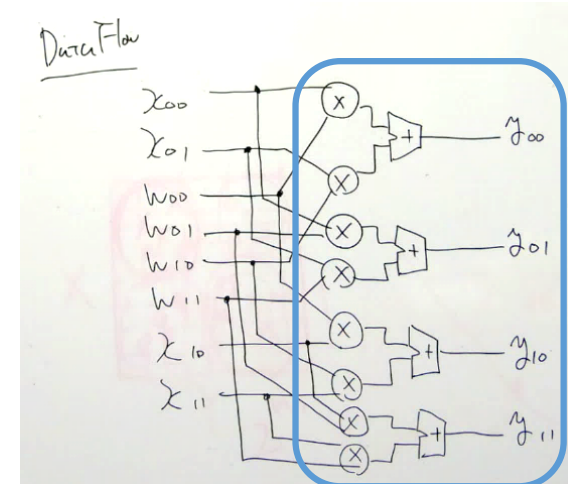
Alternatively,
 you can write a C-code

```c
for(i=0;i<2;++i){
  for(j=0;j<2;++j){
    y[i][j] = x[i][j] * w[i][j];

    // Compute terms
    for(i=0;i<2;i++){
      for(j=0;j<2;j++){
        term = 0;
        for(k=0;k<2;k++)
          term = term + x[i][k]*w[k][j];
        y[i][j] = term;
      }
    }

  }
}
```

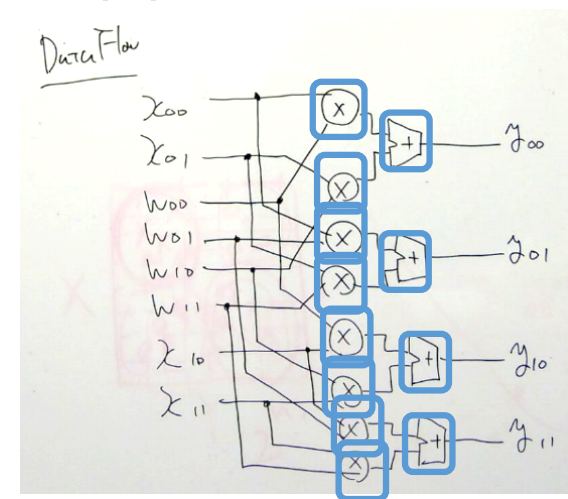# Data path description

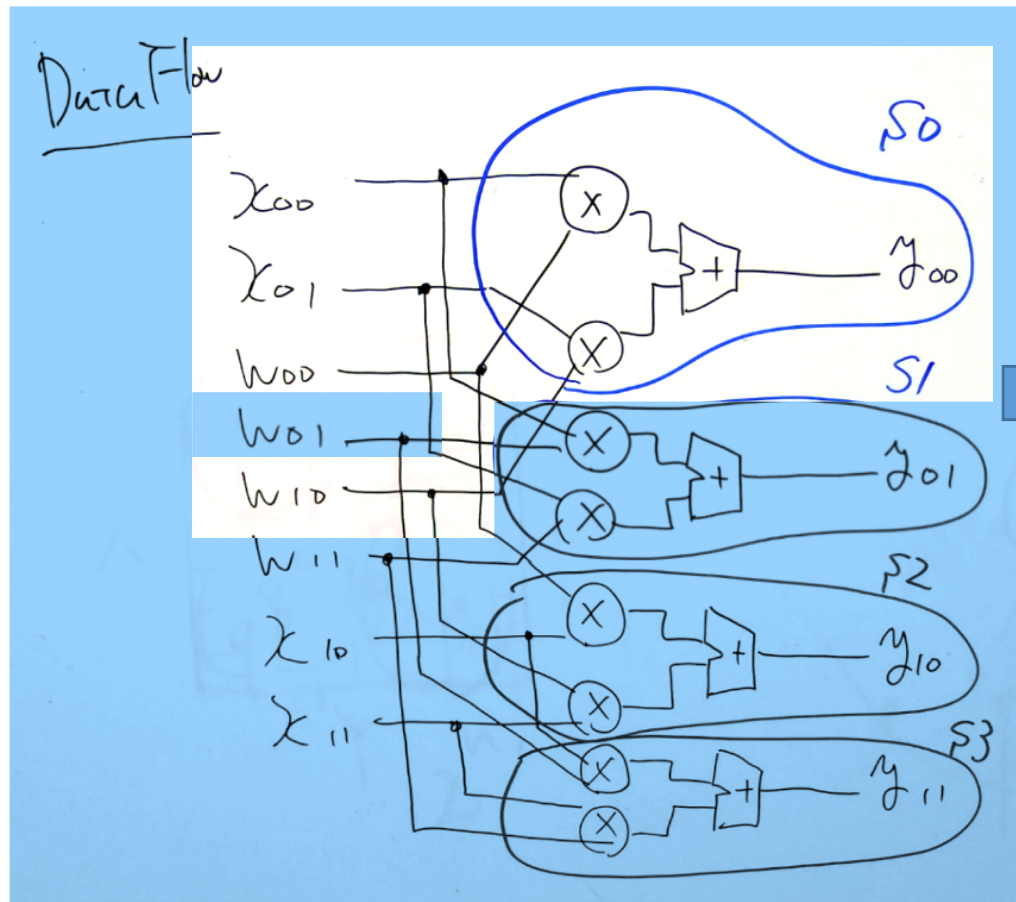# Scheduling



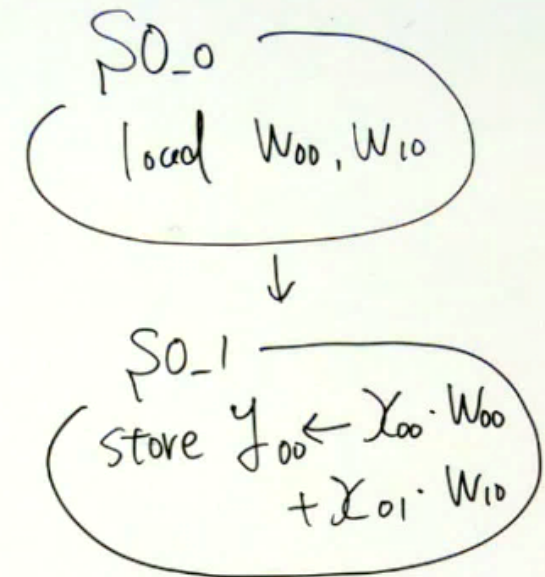MACs (Multiply ACcumulations) realization



Fully parallel realization



Sequential realization

# Finite State Machine (FSM)

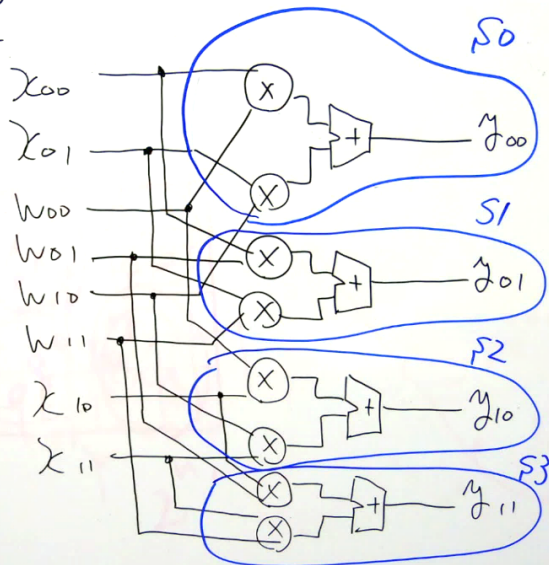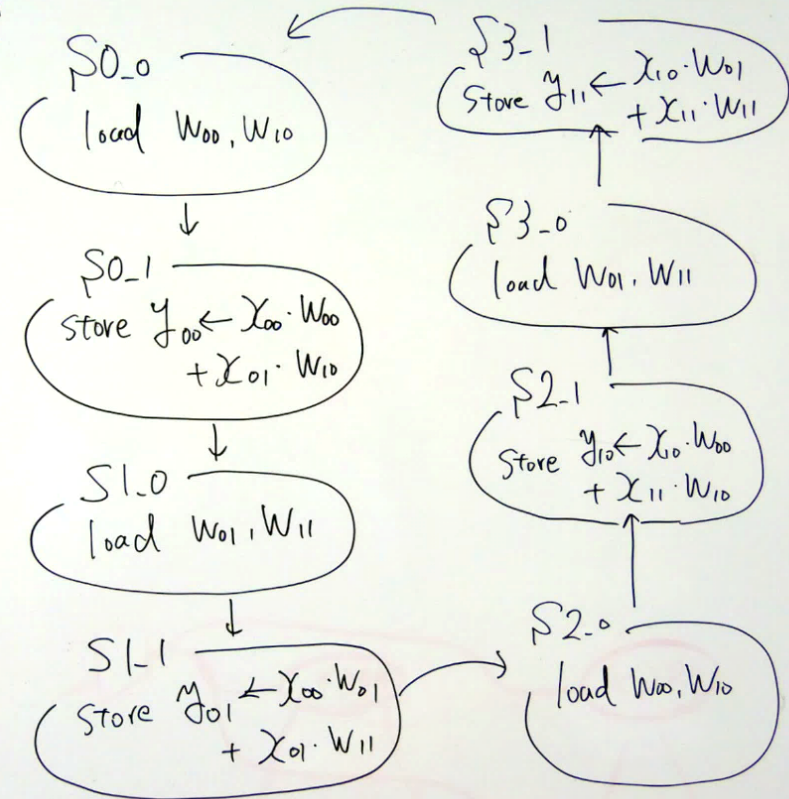# FSM (Cont.)
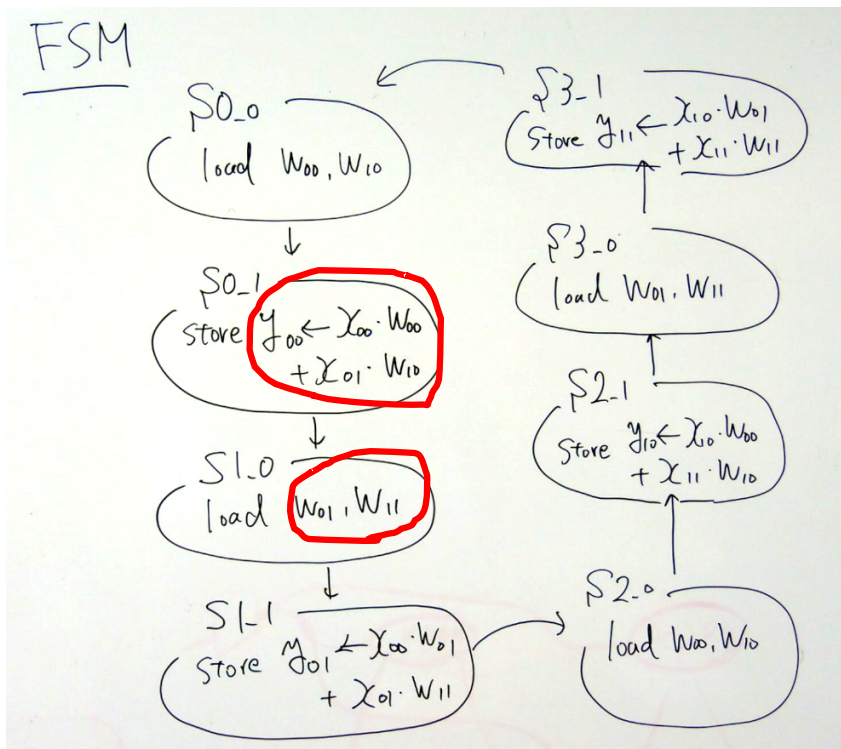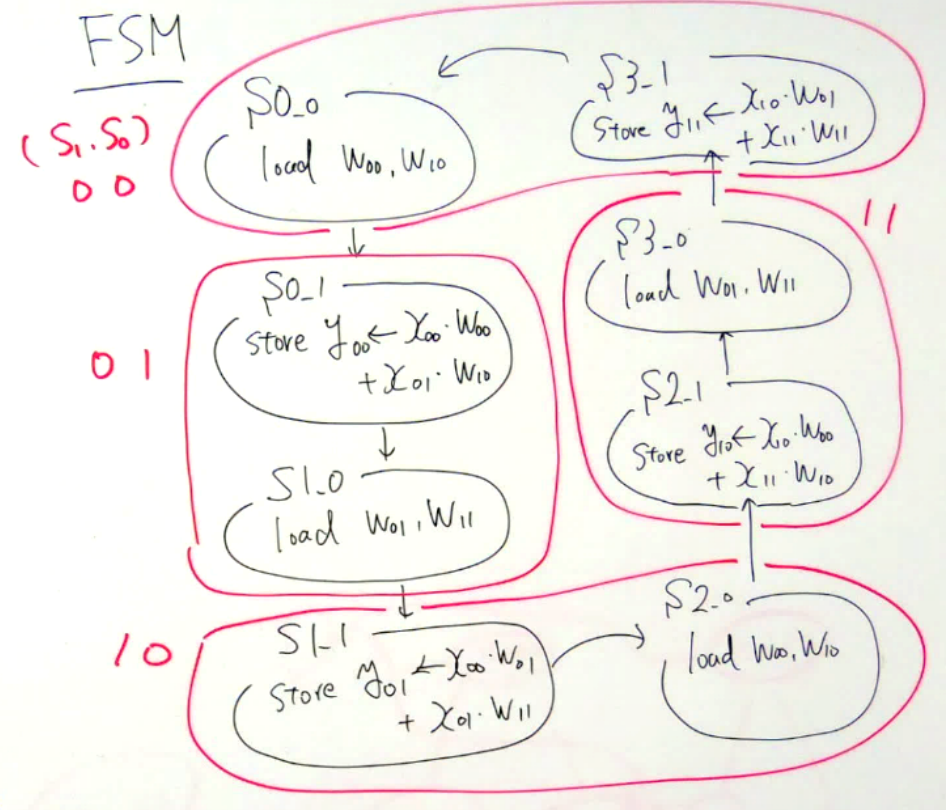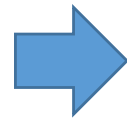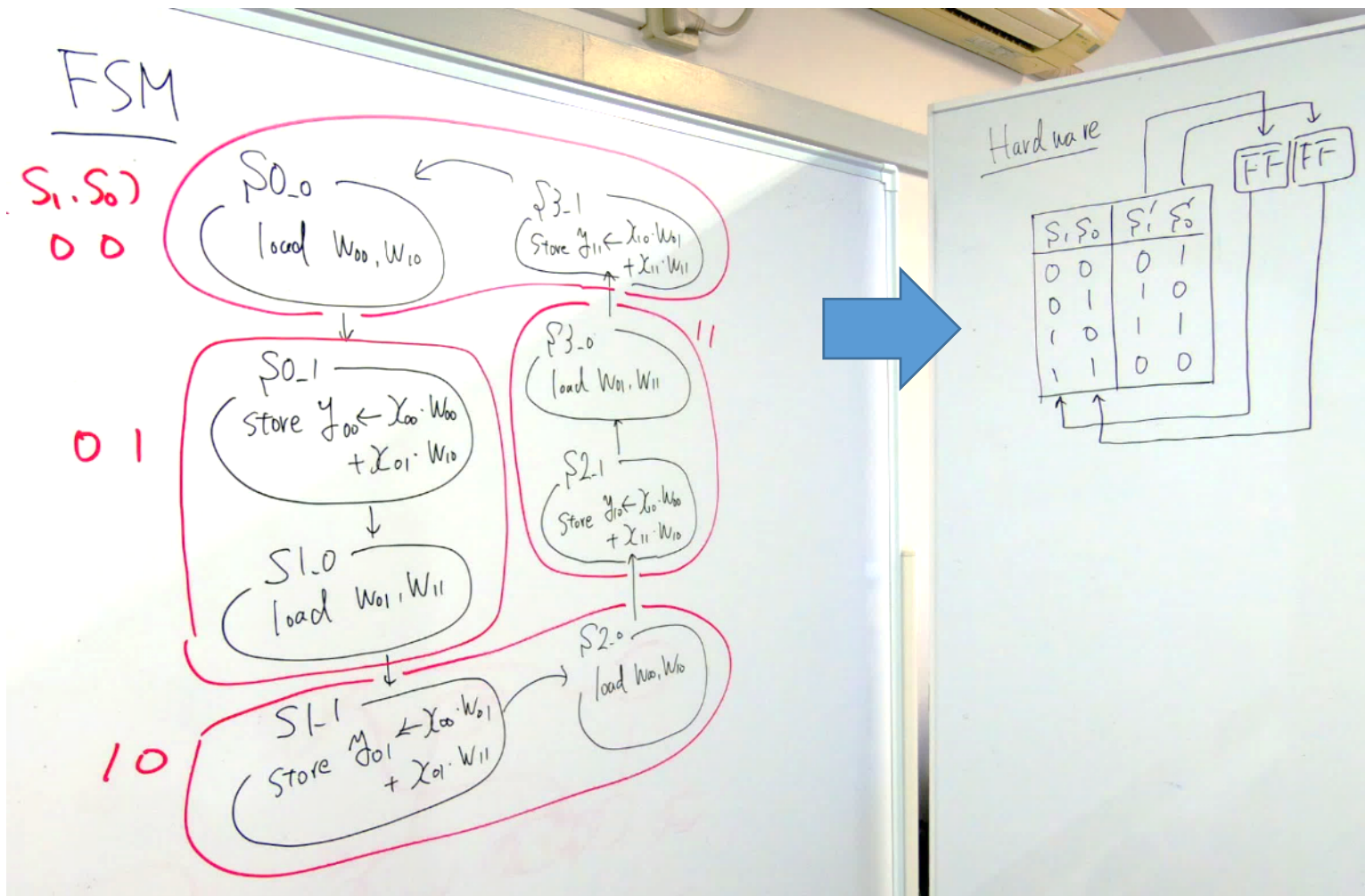
# FSM Optimization


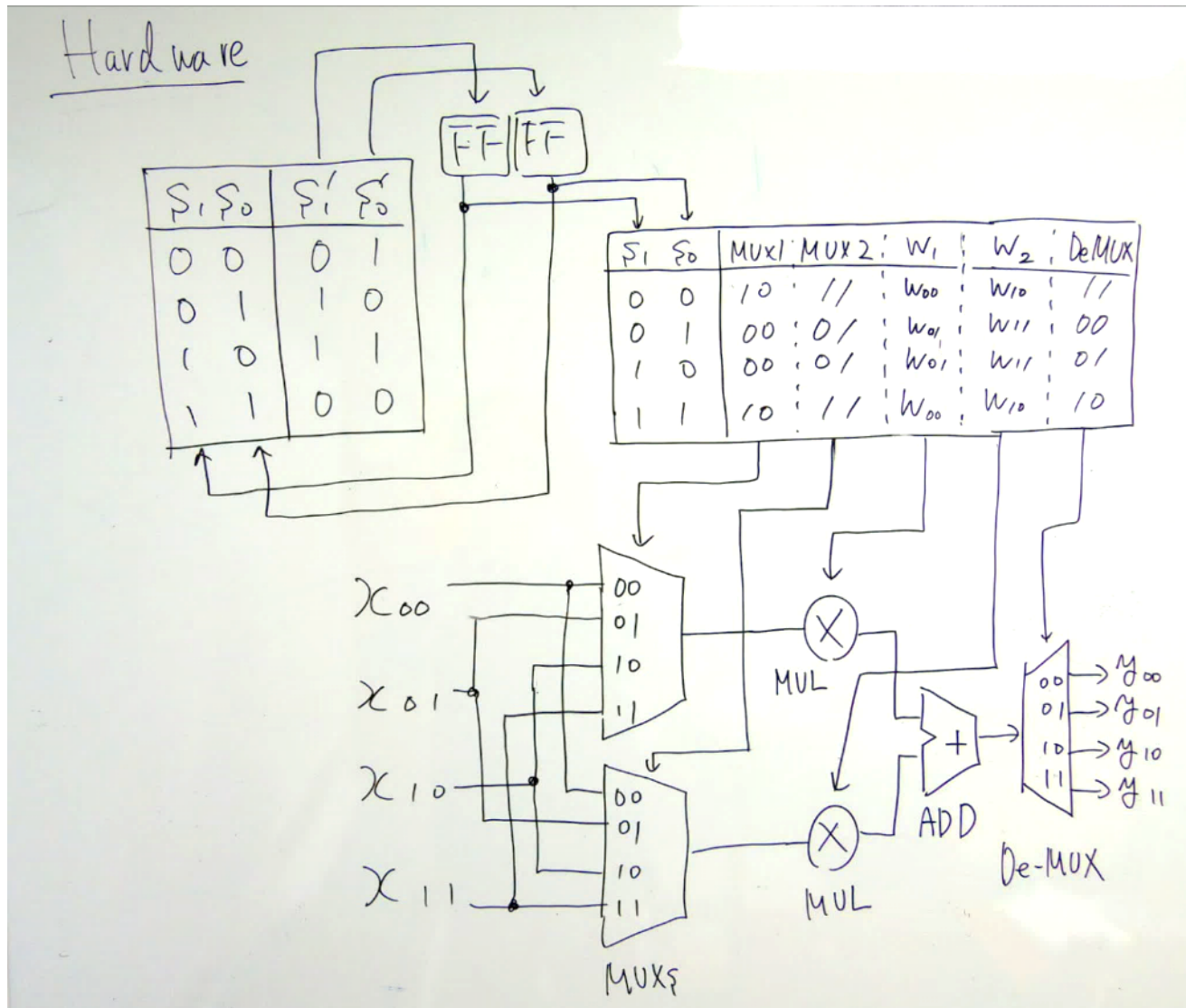
Load weights during
the MAC operation



Assign state variables
(1-hot code, gray one,
natural binary one)

# Realization of the FSM

# Assign to Primitives



Control

Data path

# Comparison of # Lines

Y = X.dot(W)+B

Python: single line!

```
for(i=0;i<2;++i){
  for(j=0;j<2;++j){
    y[i][j] = x[i][j] * w[i][j];

    // Compute terms
    for(i=0;i<2;i++){
      for(j=0;j<2;j++){
        term = 0;
        for(k=0;k<2;k++)
          term = term + x[i][k]*w[k][j];
        y[i][j] = term;
      }
    }


  }
}
```

C/C++: ten lines

```verilog
module mat_add(
    input clk, reset,
    input [7:0]x[0:3],
    output [7:0]y[0:3]
);

reg [1:0]state;
reg [1:0]mux1, mux2;
reg [7:0]w0, w1;
reg [1:0]de_mux;

always@( posedge clk or posedge rst)begin
    if( rst == 1'b1)begin
        state <= 2'b00
    end else begin
        case( state)
        2'b00:begin
            state <= 2'b01;
            mux1 <= 2'b10;
            mux2 <= 2'b11;
            w0 <= 8'b00101000;
            w1 <= 8'b11000101;
            de_mux <= 2'b11;
        end
        2'b01:begin
            state <= 2'b10;
            mux1 <= 2'b00;
            mux2 <= 2'b01;
            w0 <= 8'b00101000;
            w1 <= 8'b11000101;
            de_mux <= 2'b00;
        end
        2'b10:begin
            state <= 2'b11;
            mux1 <= 2'b00;
            mux2 <= 2'b01;
            w0 <= 8'b00101000;
            w1 <= 8'b11000101;
            de_mux <= 2'b01;
        end
        2'b11:begin
            state <= 2'b00;
            mux1 <= 2'b00;
            mux2 <= 2'b01;
            w0 <= 8'b00101000;
            w1 <= 8'b11000101;
            de_mux <= 2'b10;
        end
        endcase
    end
end

wire [15:0]mul1, mul2;
wire [16:0]w_add;

assign mul1 = w0 * mux( mux1,x[0],x[1],x[2],x[3]);
assign mul2 = w1 * mux( mux2,x[0],x[1],x[2],x[3]);

assign w_add = mul1 + mul2;

assign y[0] = (de_mux == 2'b00) ? w_add : 2'bzz;
assign y[1] = (de_mux == 2'b01) ? w_add : 2'bzz;
assign y[2] = (de_mux == 2'b10) ? w_add : 2'bzz;
assign y[3] = (de_mux == 2'b11) ? w_add : 2'bzz;

endmodule
```
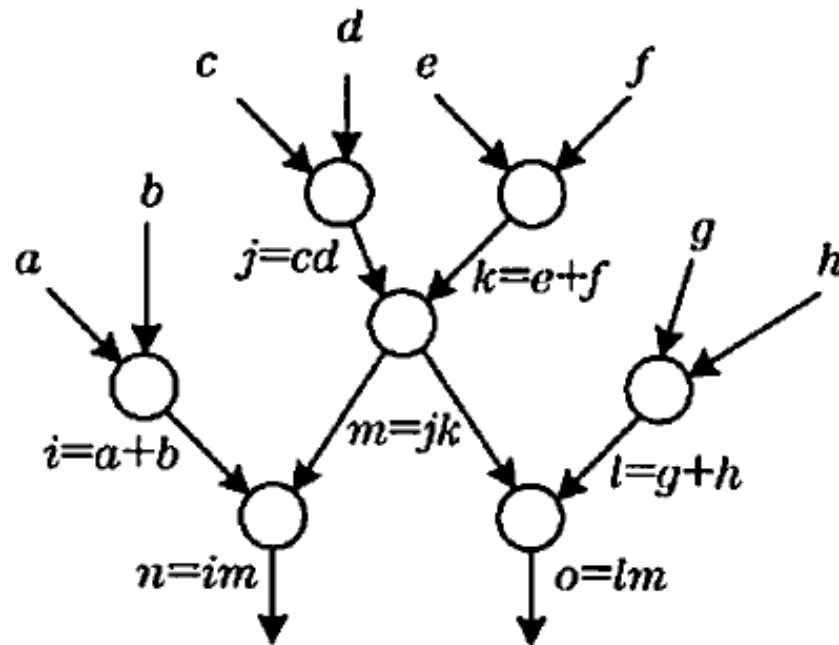
Verilog-HDL: 66 lines

# Boolean Network

- Representation of a combinational logic circuit using a directed graph without a cycle
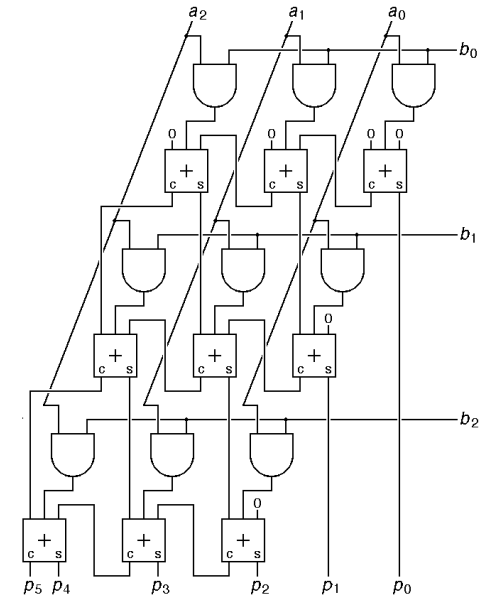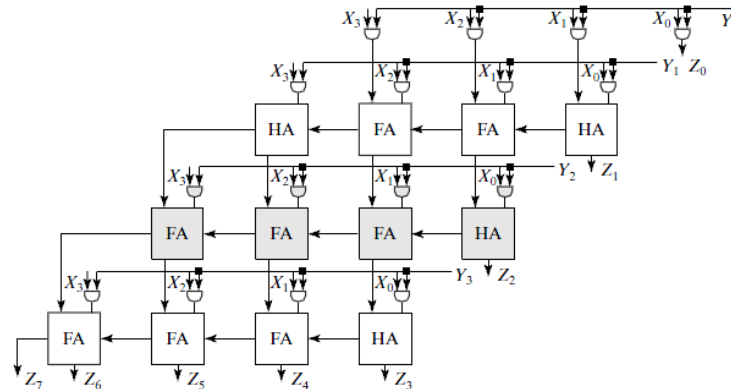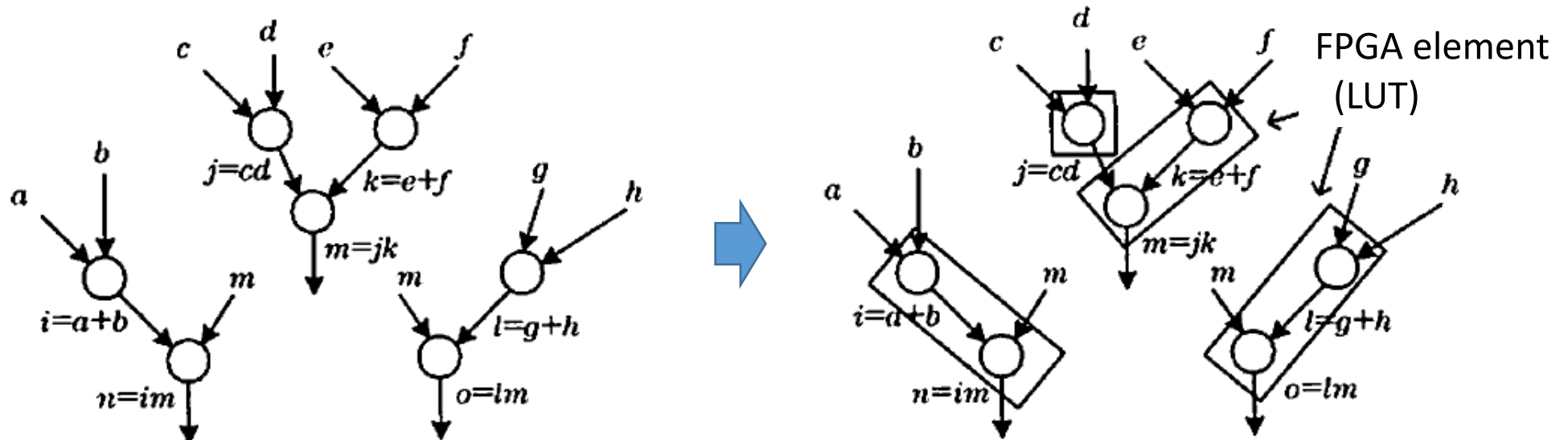- Vertex：Logic gate, Edge：Input or output

# Logic Synthesis

- Synthesize from a given HDL specification to a Boolean network
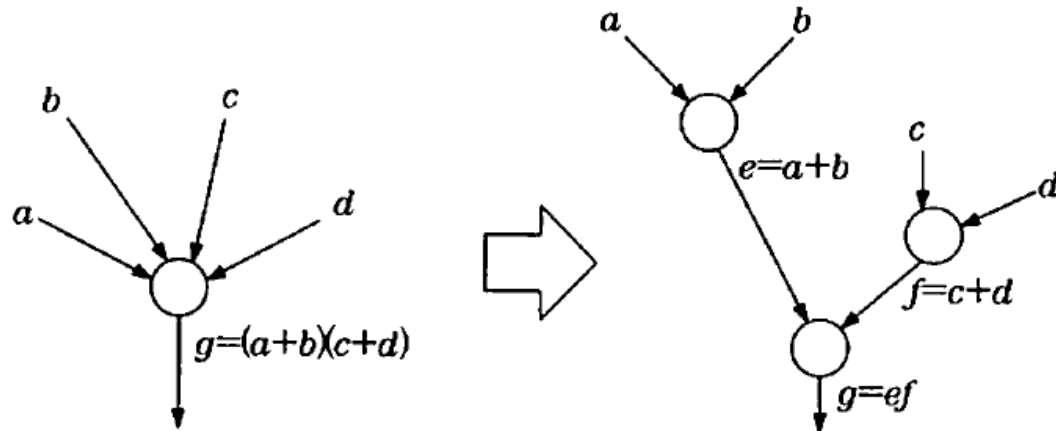
input [3:0]X,Y;
output [7:0]Z;
Z=X*Y

# Technology Mapping

- A kind of a graph covering problem
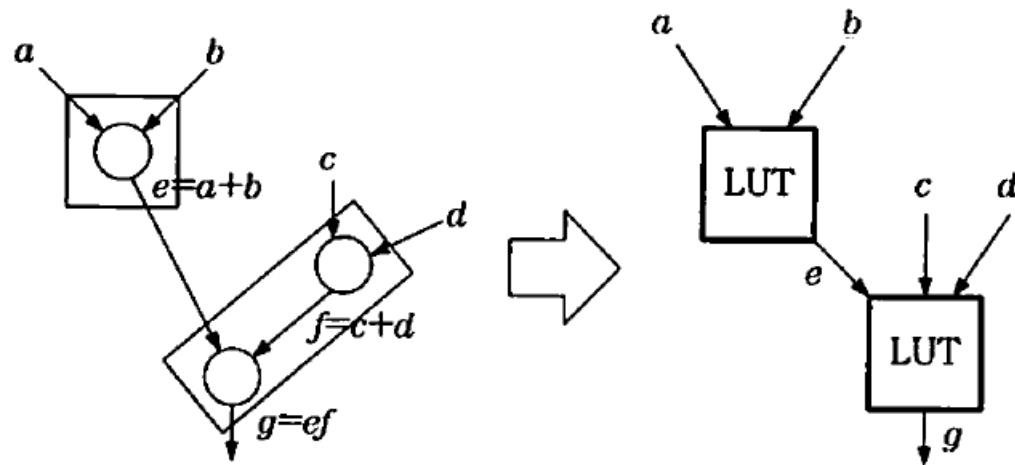- Goal: A depth optimized one by using a dynamic programming
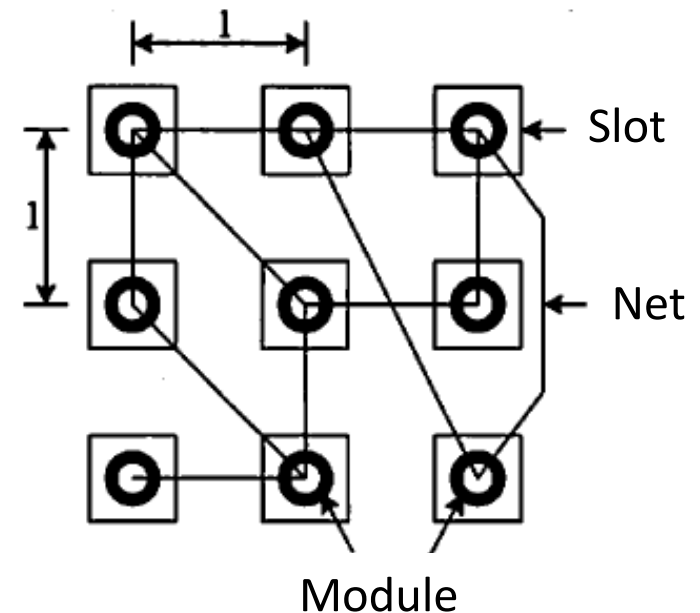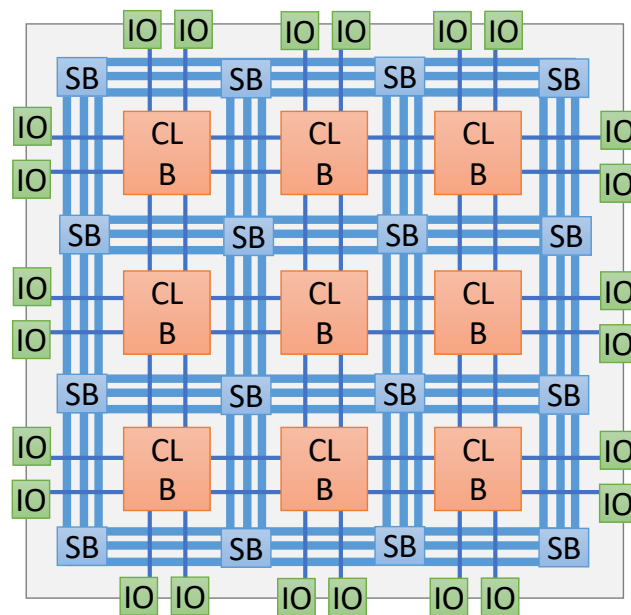
# Decomposition and Covering
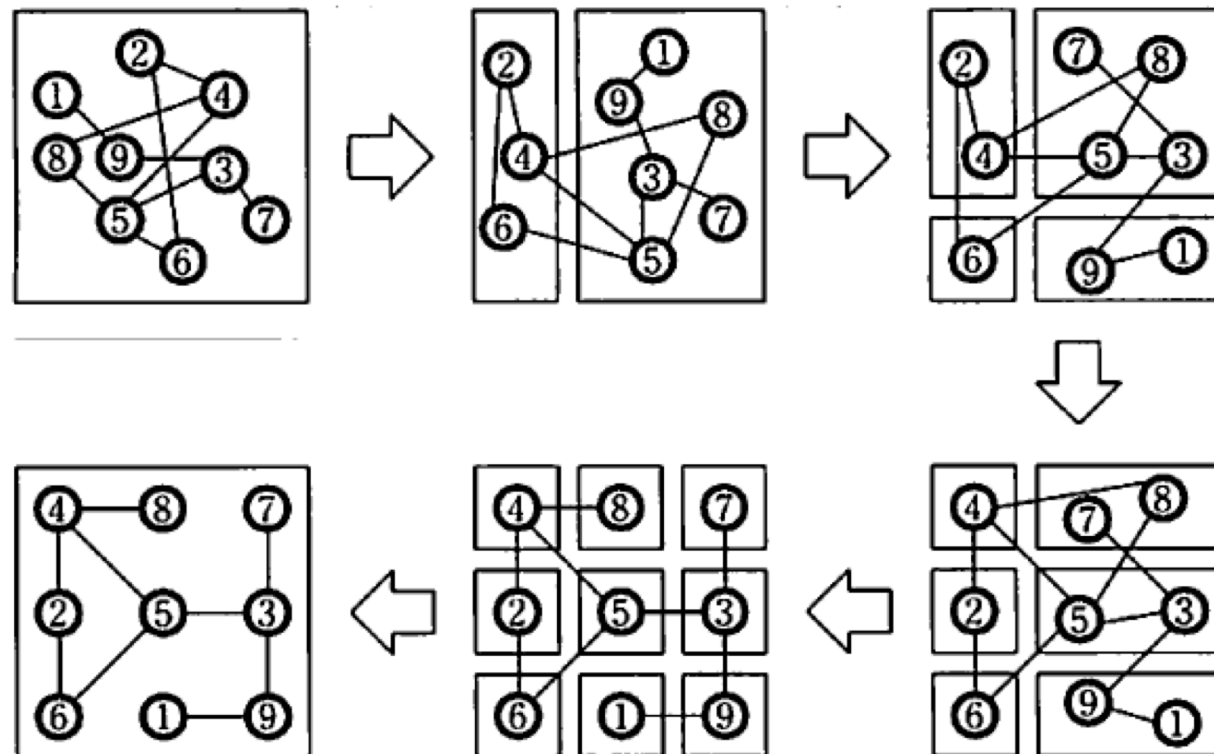
Decomposition



Covering

# Placement

- Problem to place the module (logic gate) into the slot (location)
  - 2D allocation problem → NP-complete
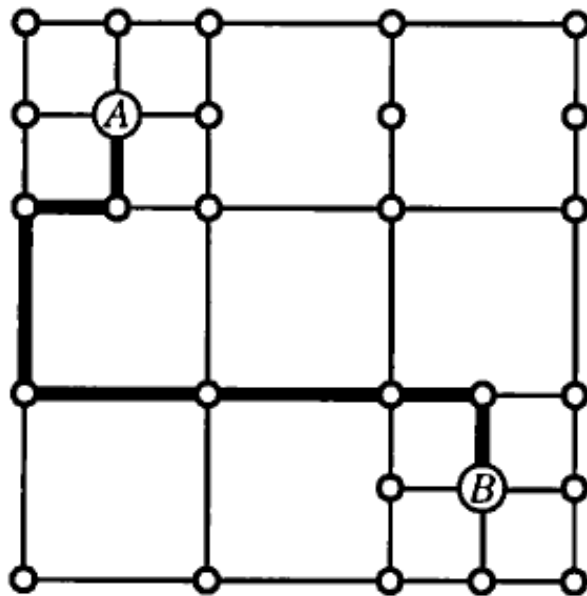  - Approximation (Simulated annealing, or min-cut tech.)

# Min-cut Placement

- Recursively divided a module set into two ones
  - Near-optimal solution can be obtained in the short time
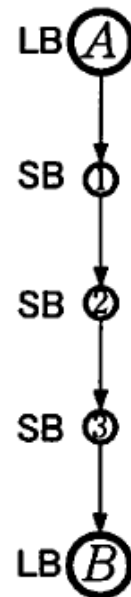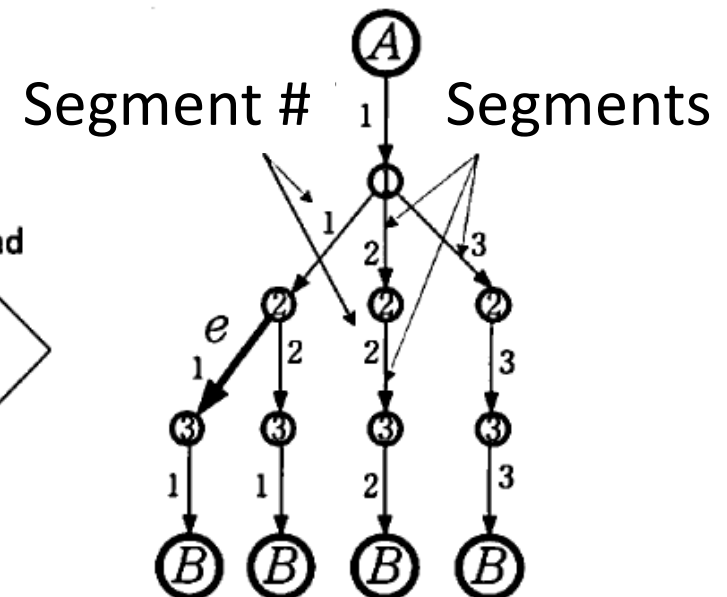- Used in a commercial FPGA design tool

# Routing

- Global routing: Determine the rough wiring path
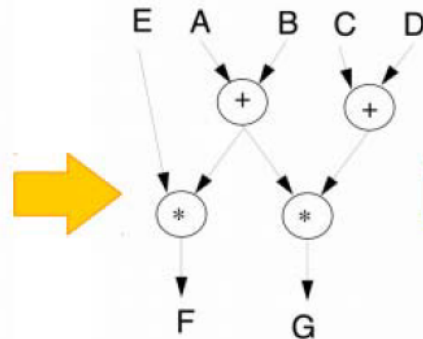- Local one: Determine the wiring segment and switch



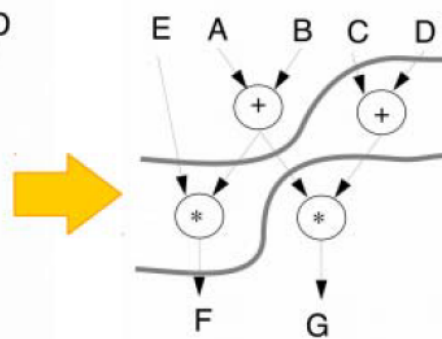Global routing

Local routing

# High-Level Synthesis (HLS)



Input Behavioral Spec.  Dataflow  Scheduling

Data-path generation

Mapping to resources (Binding)

Controller (FSM) Generation

# System on Chip FPGA

- Xilinx: Zynq, Intel: Intel SoC



Data Storage by DDR Memory

USB macro

Re-use software libraries by operating system

Source: Xilinx Inc. Zynq-7000 All Programmable SoC

# Conventional Design Flow for the SoC FPGA



1. Behavior design

2. Profile analysis

3. IP core generation by HLS

4. Bitstream generation by FPGA CAD tool

5. Middle ware generation

Source: Xilinx Inc.

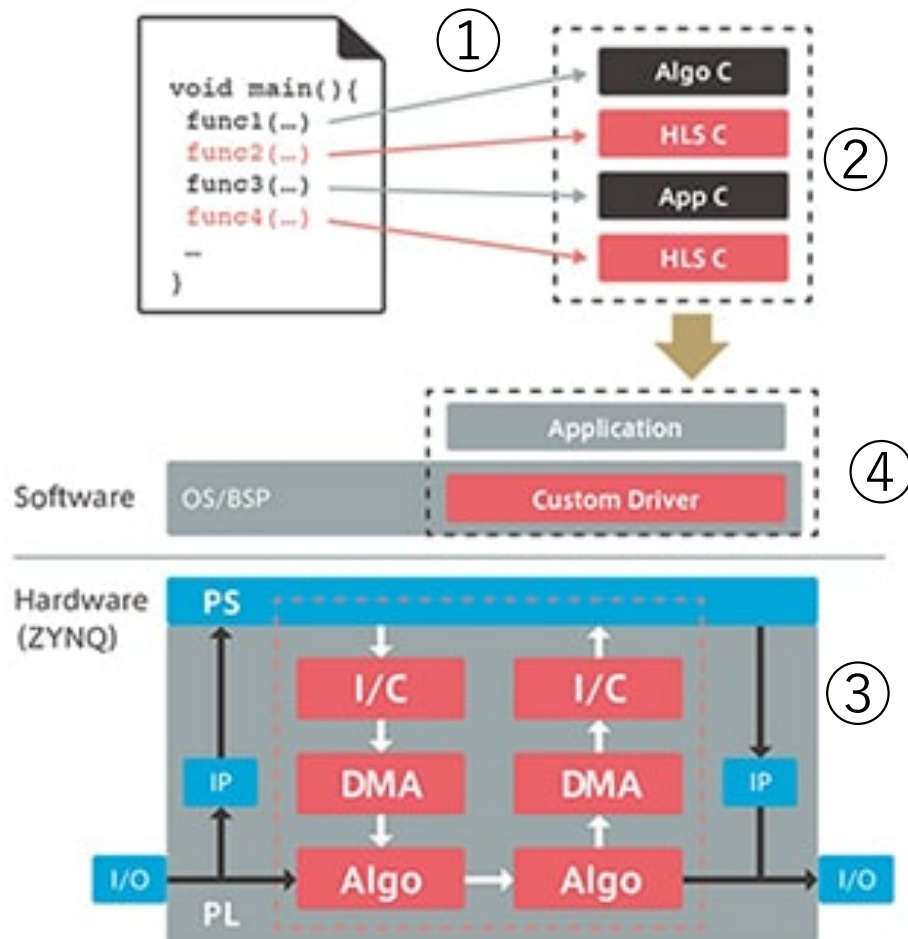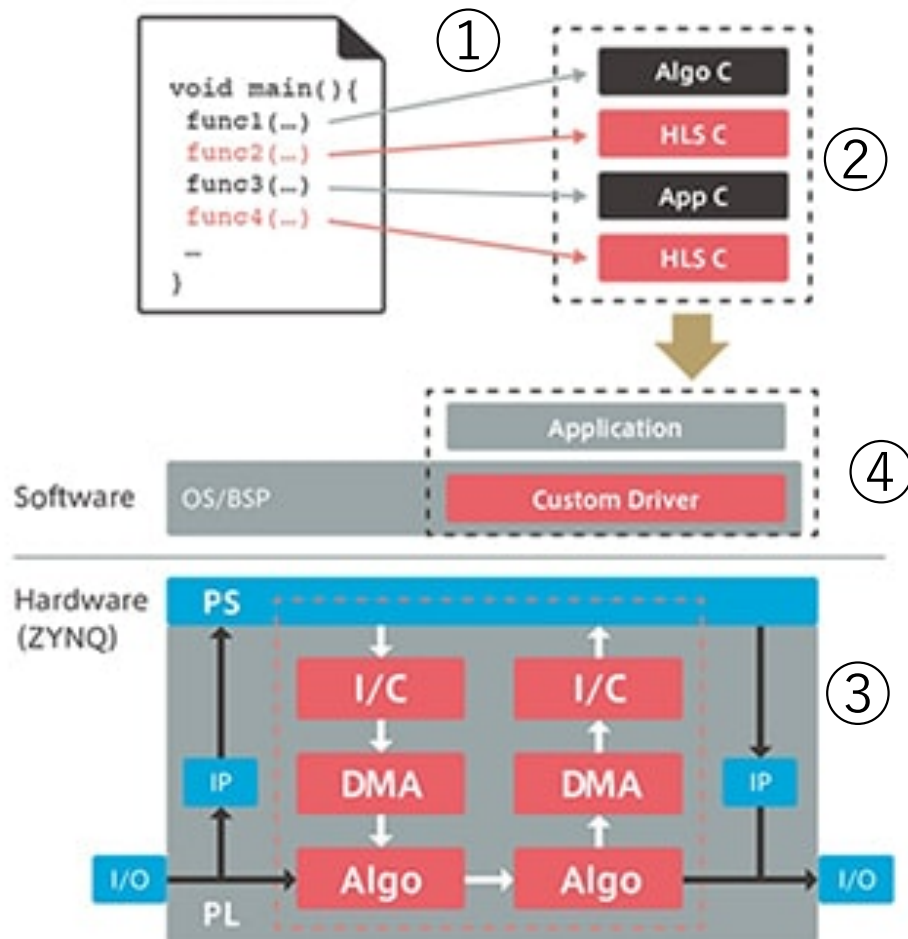# System Design Tool for the SoC FPGA



1. Behavior design

    **+ pragmas**

2. Profile analysis

3. IP core generation by HLS

4. Bitstream generation by
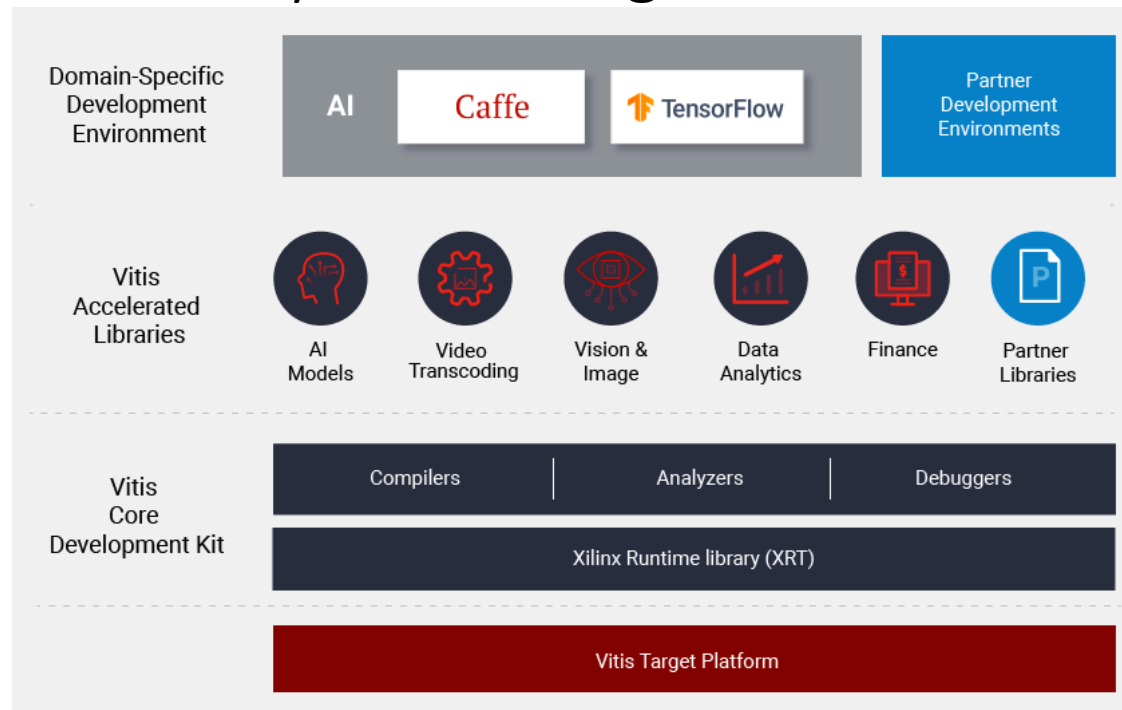    FPGA CAD tool

5. Middle ware generation

    ↓

    Automatically done

Source: Xilinx Inc.

# Design Flow for Versal Architecture:
 e.g. AI Accelerator

- Application specific framework (Caffe, TensorFlow)

- Convert pre-trained model with Vitis Libraries

- Conventional system-design tool-flow



Source: Xilinx Inc.

# Summary

- FPGA: Reconfigurable LSI or Programmable Hardware

- It consists of a programmable logic array

   and a programmable interconnection

- Standard FPGA design supports an RTL based one

- Benefits: Productivity, lower non-recurring engineering costs, maintainability, faster time to market

# Homework 1

1. (Mandatory) Show an application using a FPGA

2. (Selective) The modern FPGA has several configuration devices as follows:
   - Static RAM (SRAM)
   - Flash memory
   - Anti-fuse

   In that case, both the SRAM-based FPGA and Flash-based one supports "multiple-time" configurations, while an anti-fuse only "one-time" configuration. Actually, some vender adopts the anti-fuse-based FPGA than SRAM-based one. Why?

Deadline is 21$^{st}$, Nov., 2019

(At the beginning of the next lecture)

Send a PDF file to nakahara{at}ict.e.titech.ac.jp