
Reconfigurable Architectures

FPGA as an accelerator

AMANO, Hideharu

hunga@am.ics.keio.ac.jp

PLD (Programmable Logic Device)

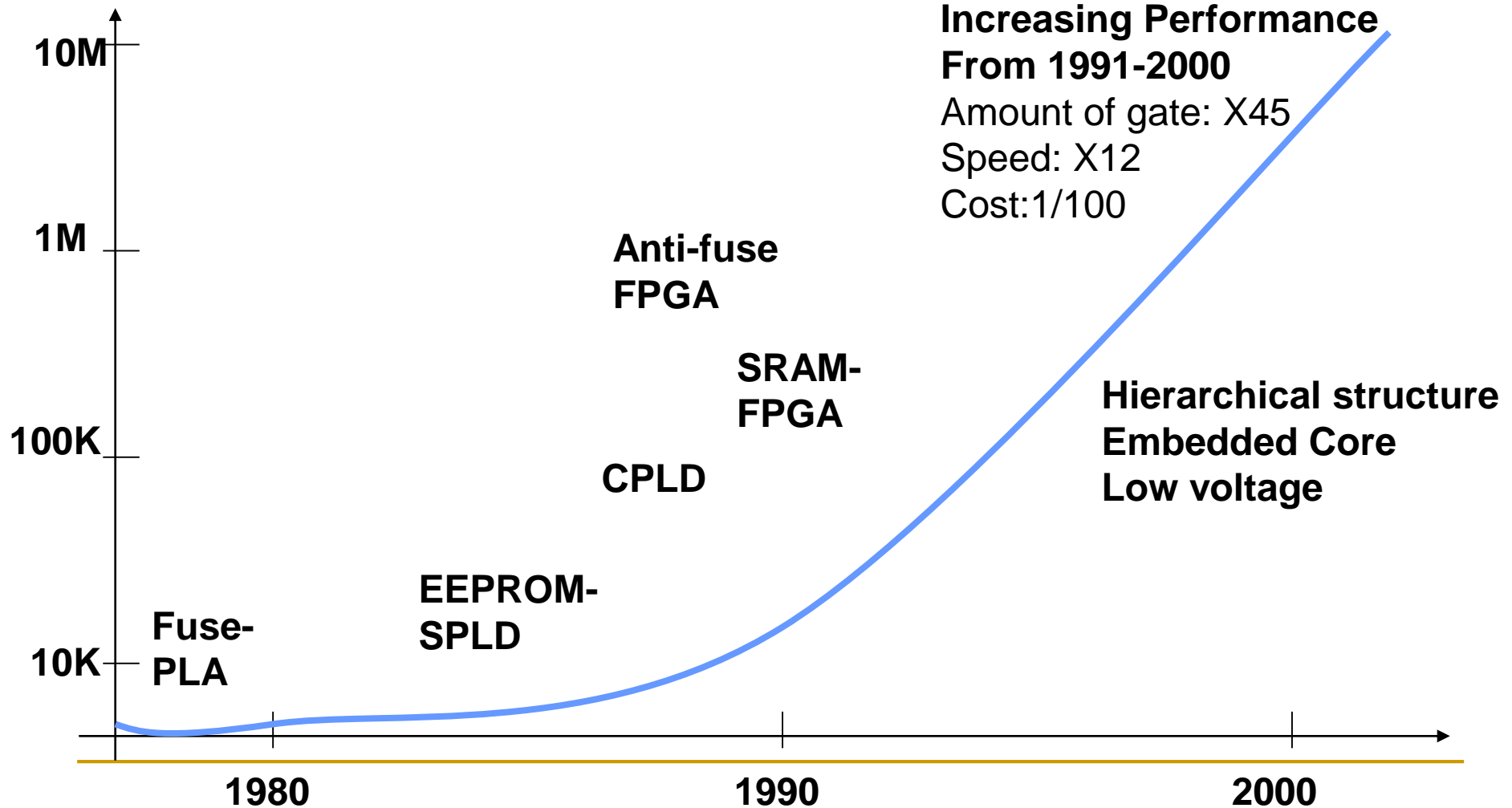
- Integrated Circuit whose logic function can be defined by users.
 - ↔ Standard IC, ASIC (Application Specific IC)
- SPLD (Simple PLD) / PLA (Programmable Logic Array)
 - Small scale IC with AND-OR array
- CPLD (Complex PLD)
 - Middle scale IC with AND-OR array
- FPGA (Field Programmable Gate Array)
 - Large scale IC with LUT



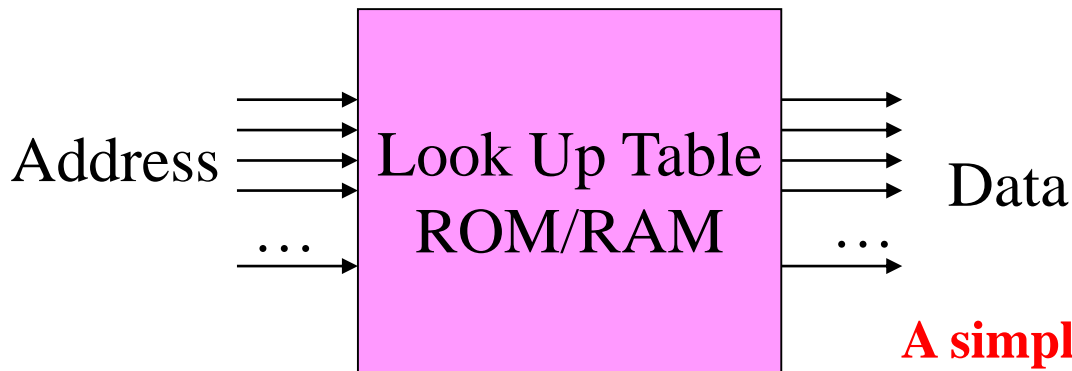
Caution! Terms are not well defined!

Rapidly development of PLD

Gate number



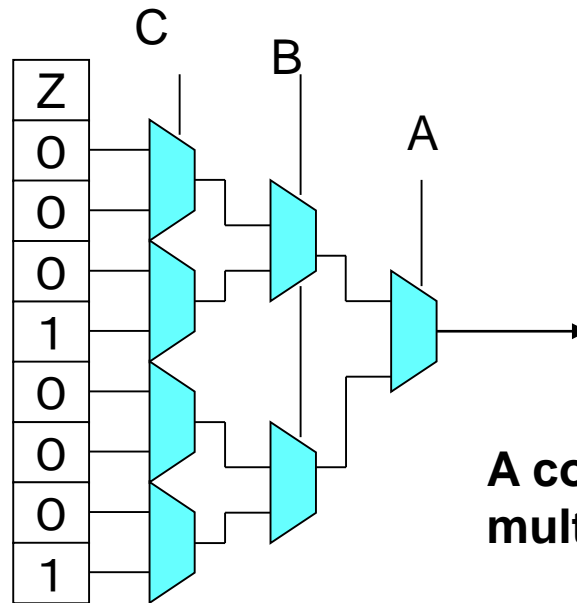
LUT: Look Up Table



A simple ROM/RAM can be used as a random logic.



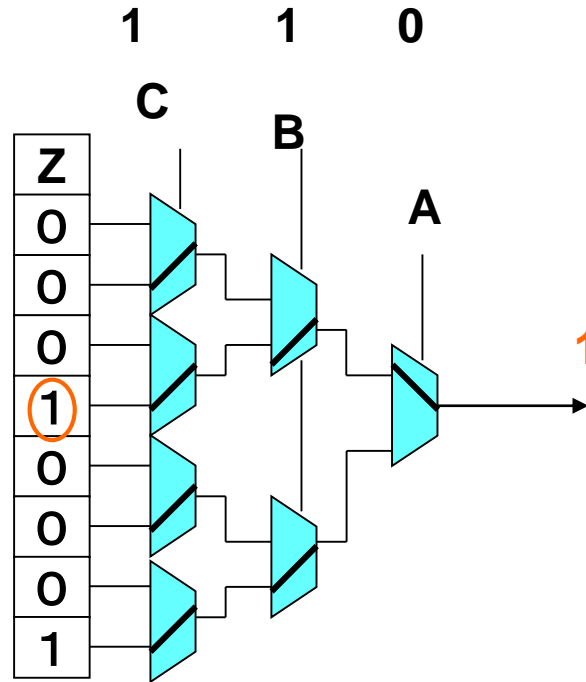
ABC	Z
000	0
001	0
010	0
011	1
100	0
101	0
110	0
111	1



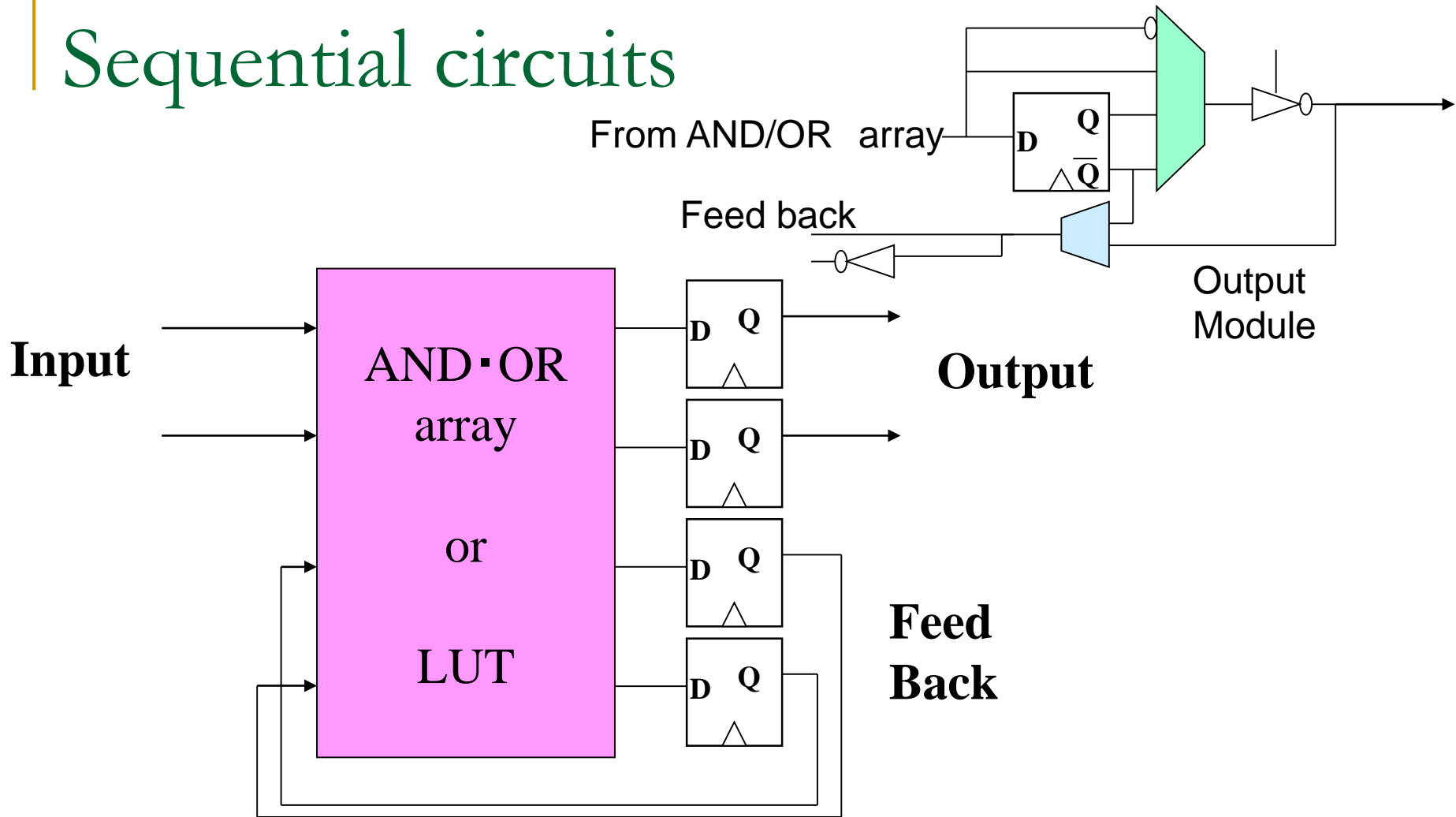
A combination of memory and multiplexers are commonly used.

An example using LUT: Look Up Table

ABC	Z
000	0
001	0
010	0
011	1
100	0
101	0
110	0
111	1

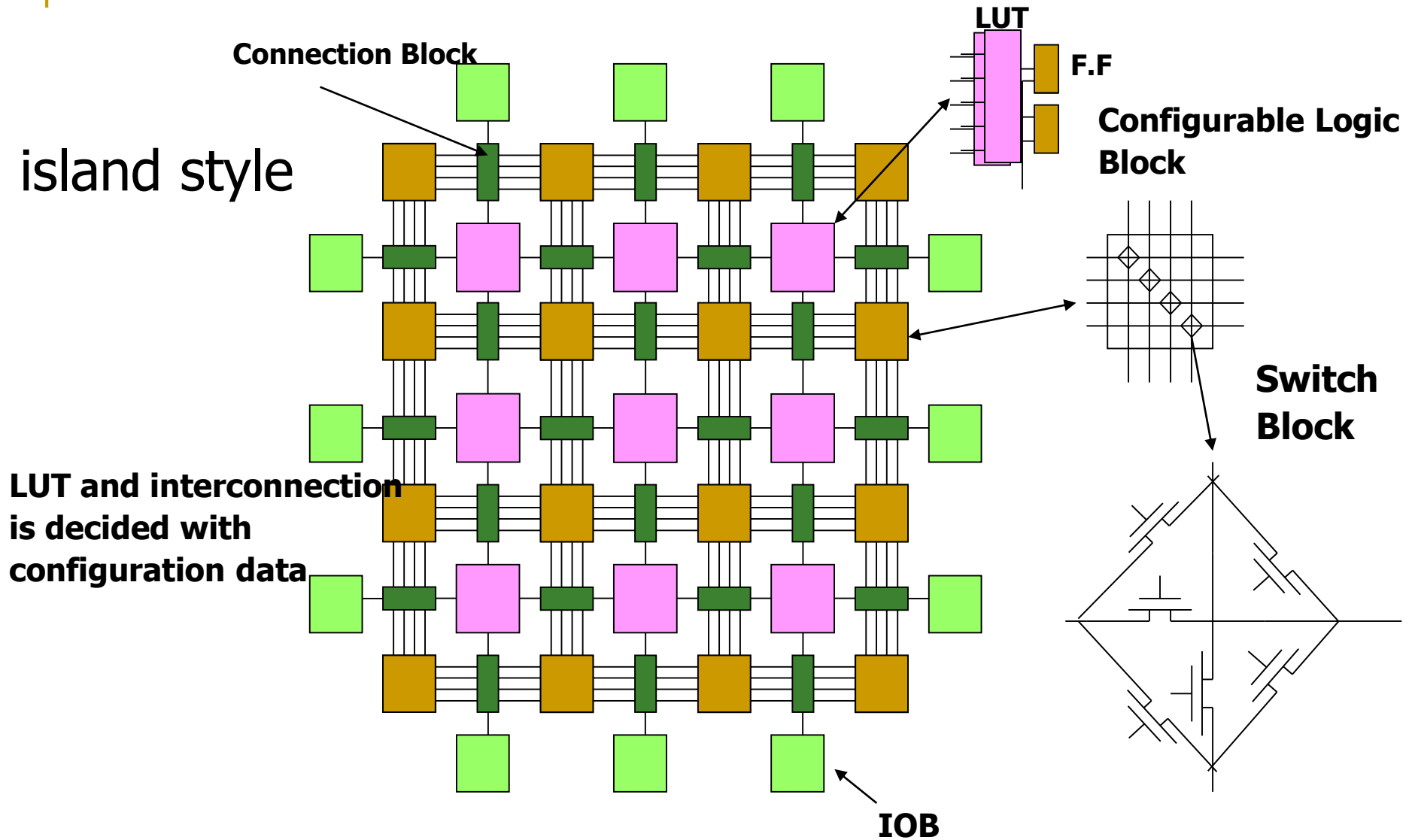


Sequential circuits

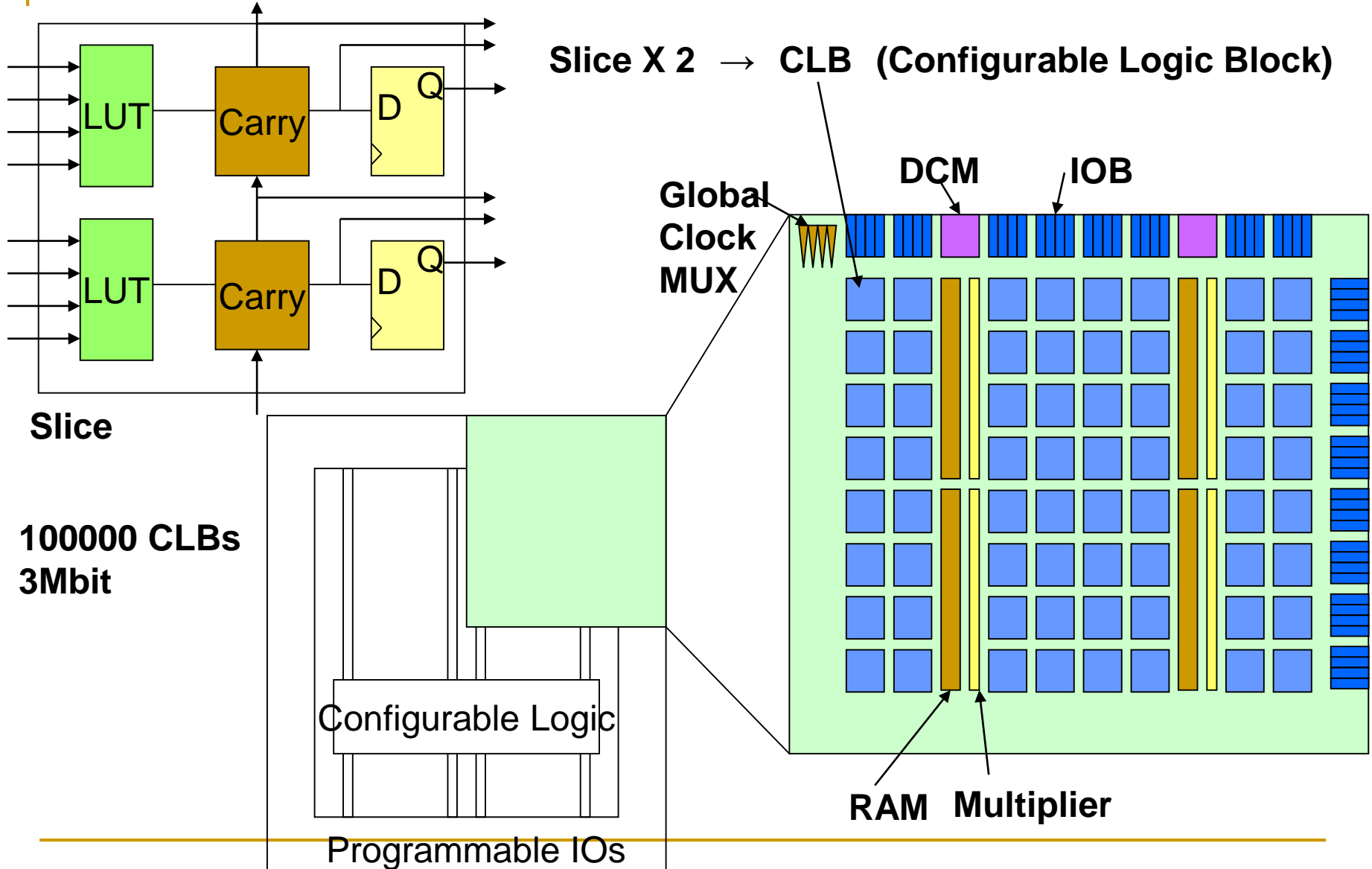


Sequential circuit (state machine) can be built by attaching Flip-flops and feed back loops.

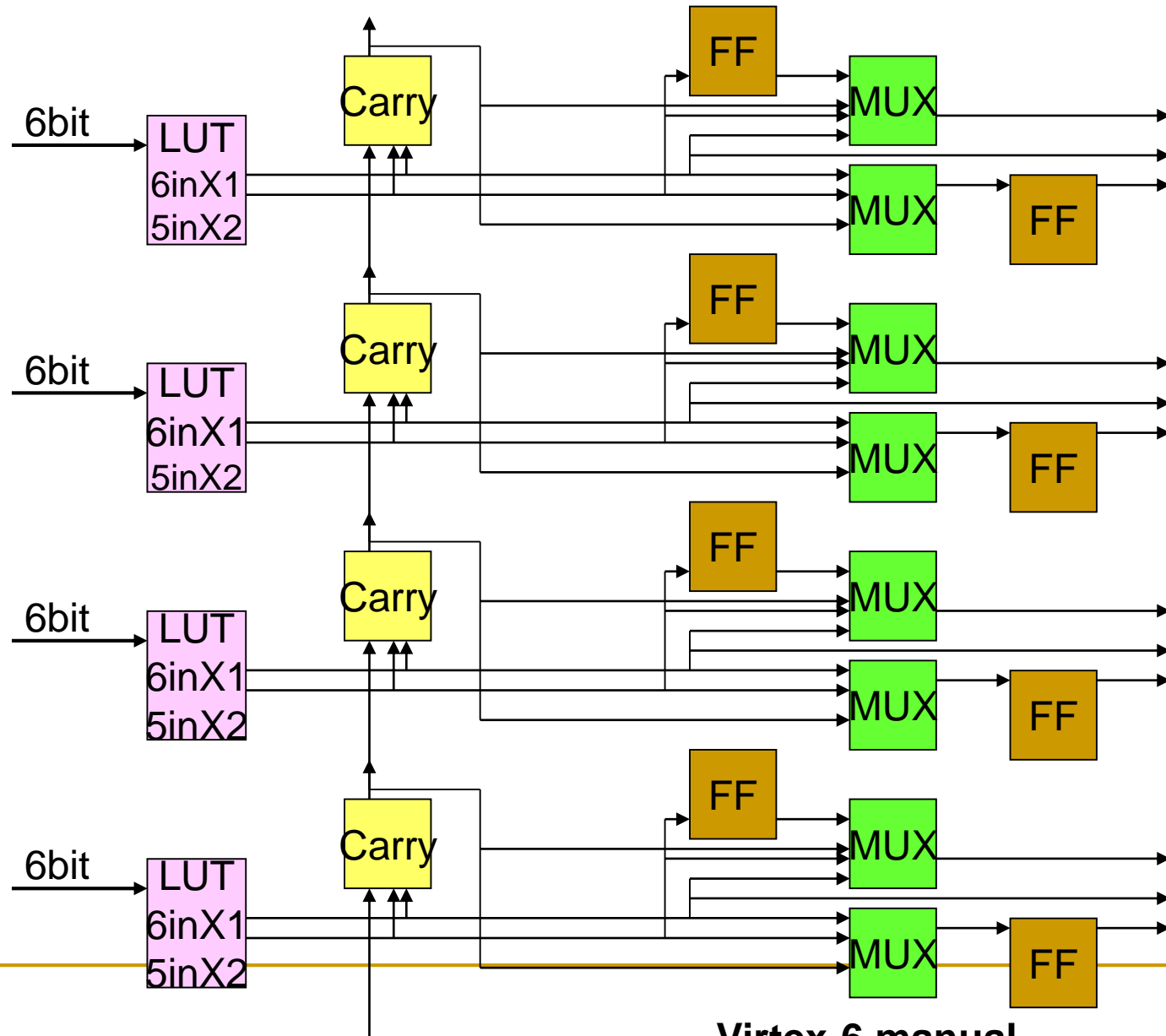
FPGA(Field Programmable Gate Array)



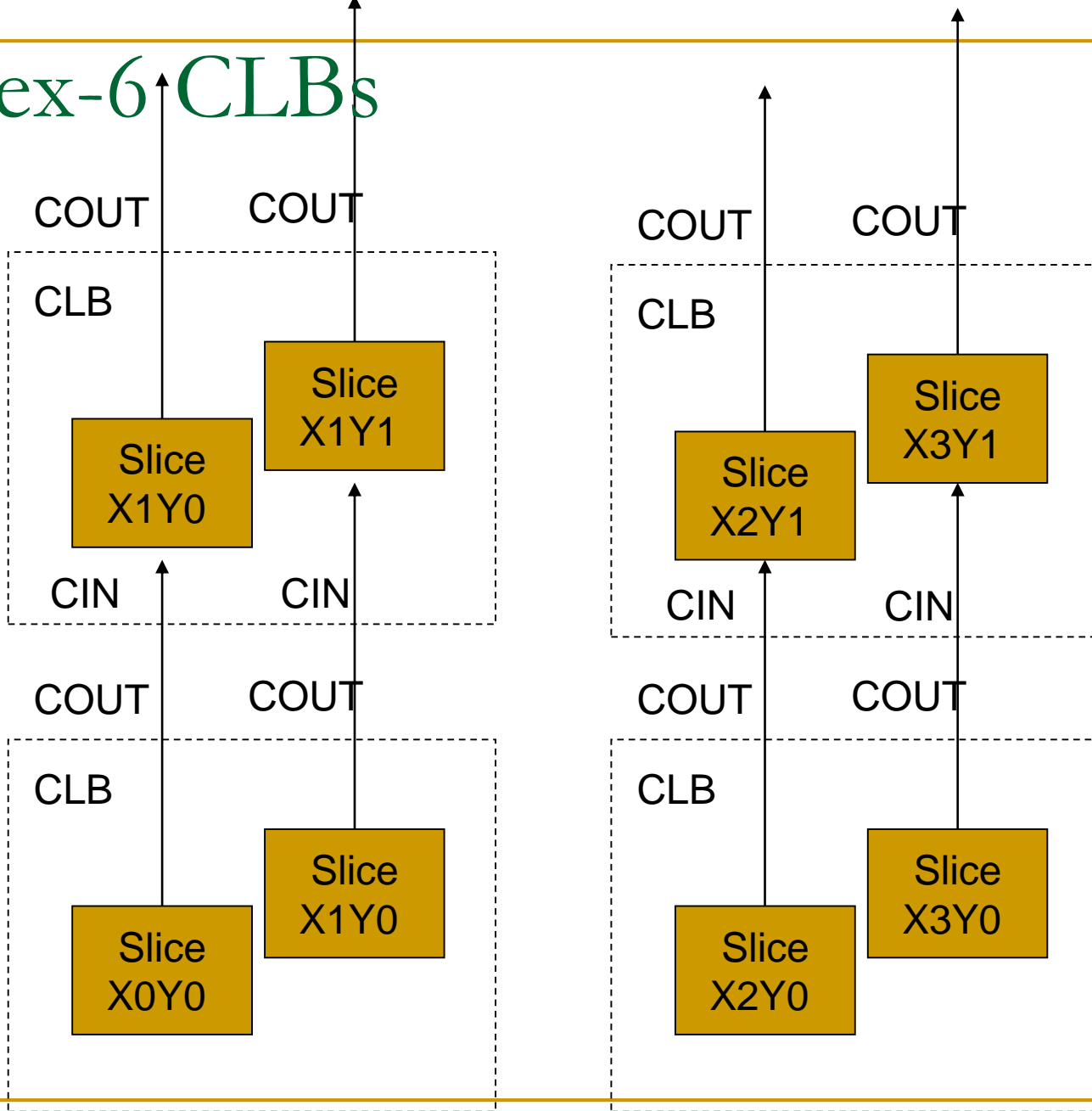
Xilinx Virtex II



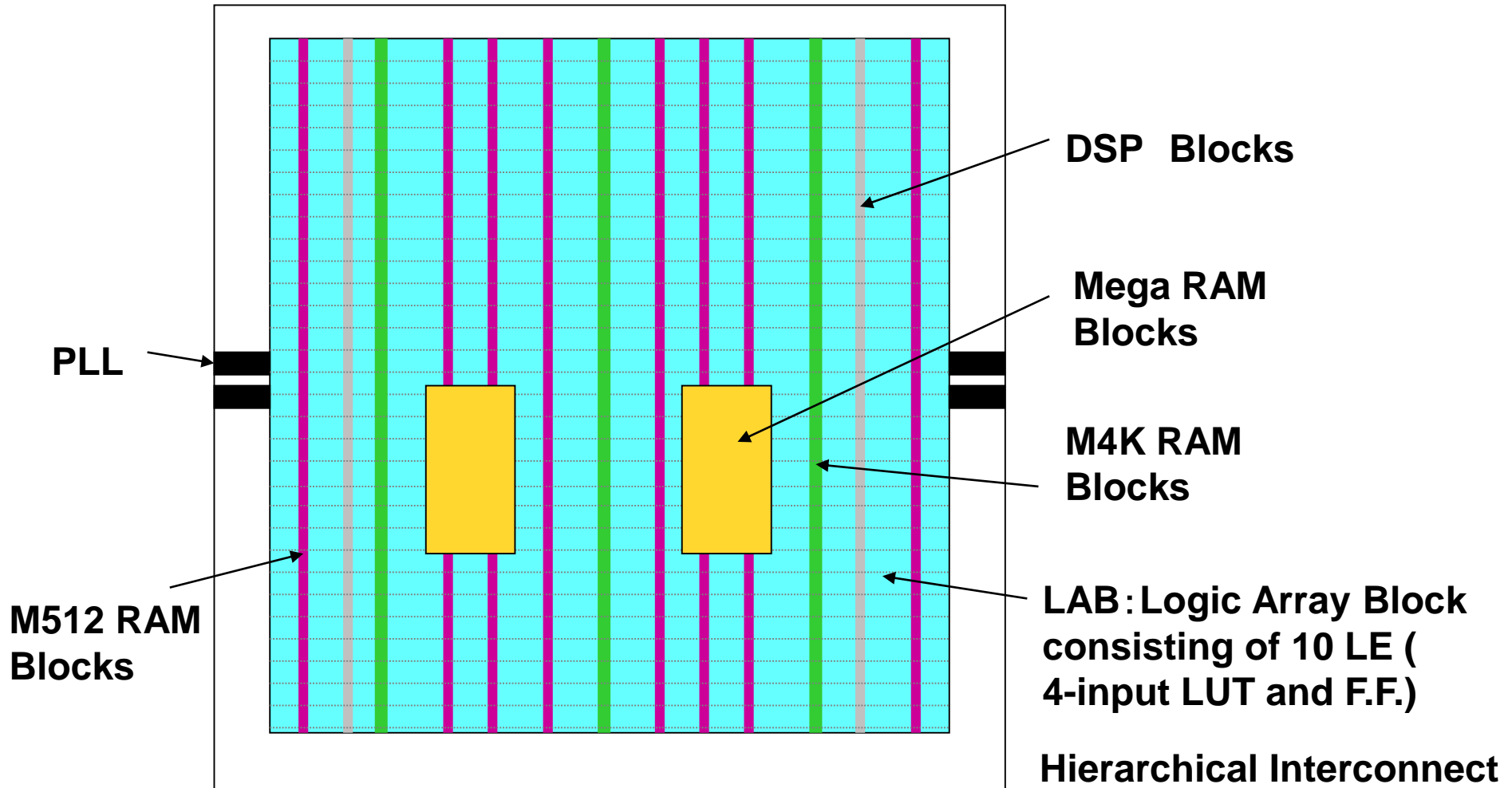
Slice structure of Virtex-6



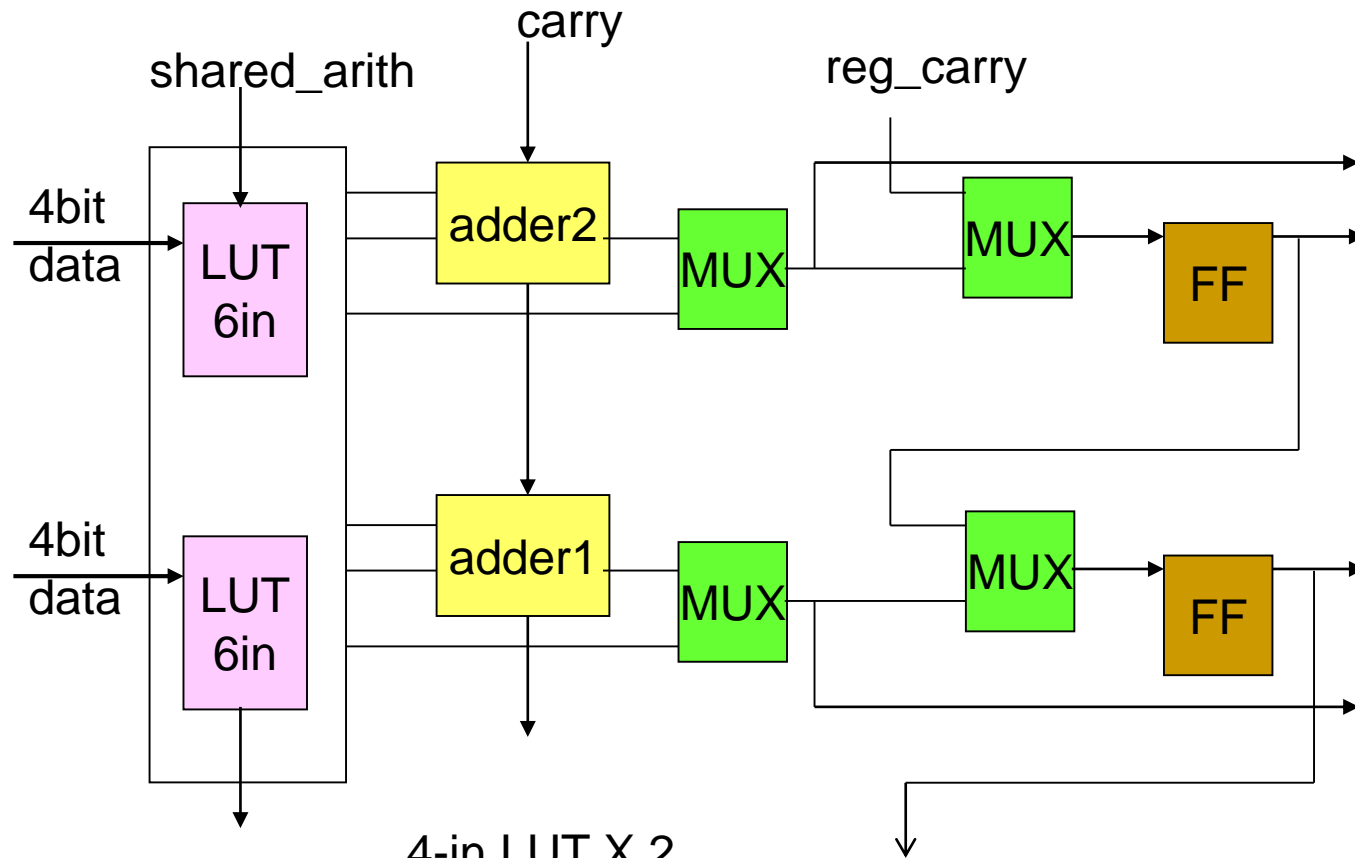
Virtex-6 CLB



Altera Stratix II



Stratix-IV ALM Structure



4-in LUT X 2

5-in LUT + 3-in LUT

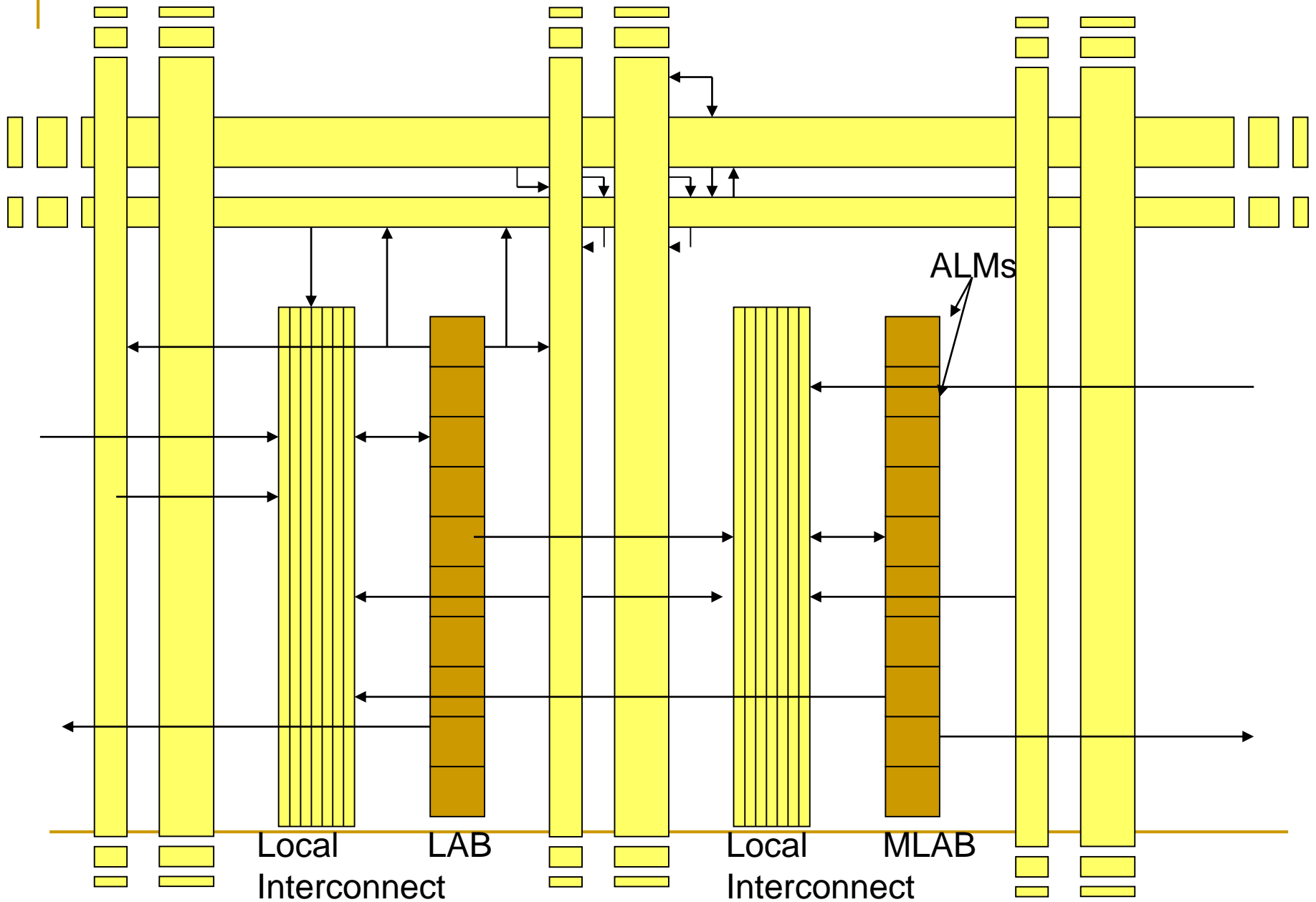
5-in LUT + 4-in LUT 1-input shared

5-in LUT + 5-in LUT 2-input shared

6-in LUT

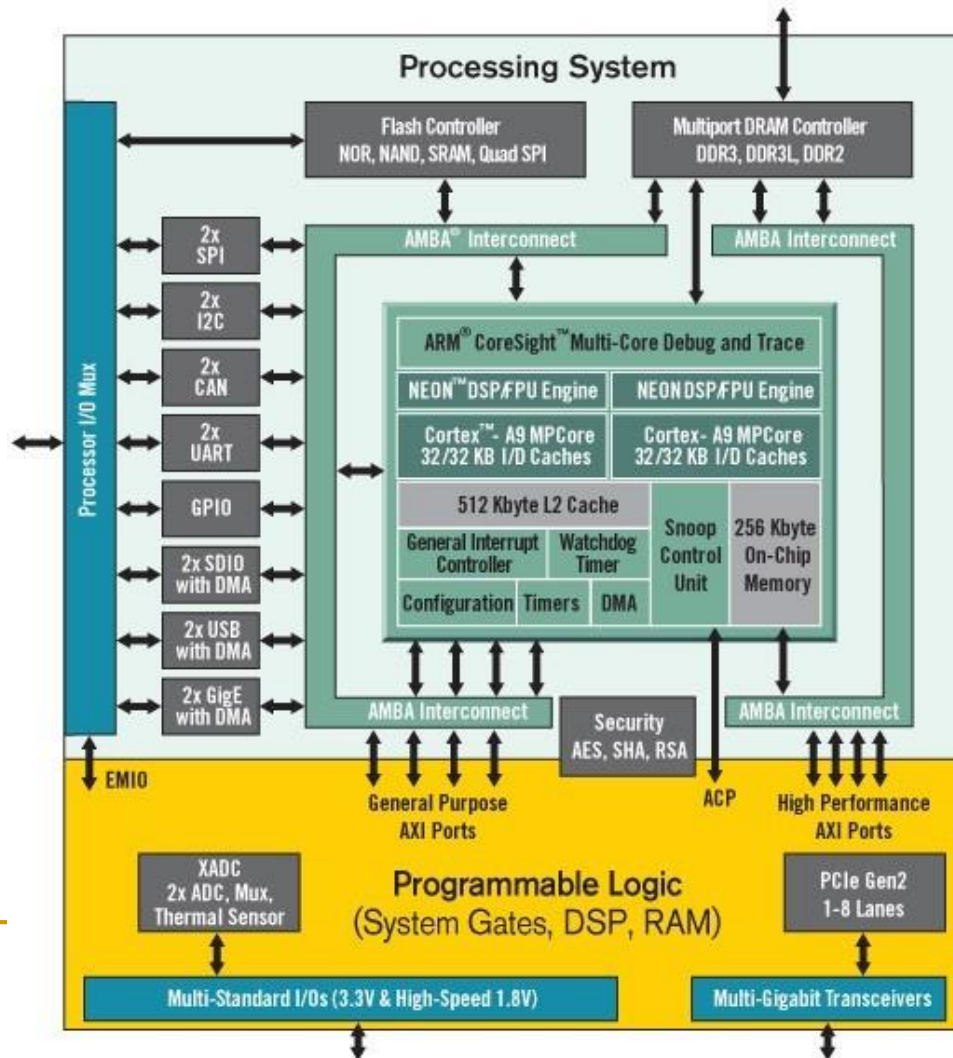
6-in LUT + 6-in LUT 4-input shared

Stratix-IV LAB structure



Zynq All programmable SoC

- ARM Cortex-A9 (PS part) Dual Core CPU +
- 28nm Artix-7/Kintex-7 based FPGA (PL part)



Technologies vs. Product

90nm

65nm 60nm

45nm 40nm

28nm

High-end
Virtex-4LX/FX/SX
200000LC

Virtex-5LX/LXT/SXT/
FXT/TXT
330000LC

Virtex-6LXT/SXT/
HXT/CXT
760000LC

Virtex-7
T/XT/HT
2000000LC

Stratix-II/GX
179400LE

Stratix-III/L/E
338000LE

Stratix-IV
/E/GX/GT
531200LE

Stratix-V
/E/GX/GS/GT
359200ALM

X1.5-X2.5 / generation

Middle range

Kintex-7
480000LC

Arria

Arria-II

Arria-IV
174000LE

Low-cost
Spartan-3A N/DSP
53000LC

Spartan-6LX/LXT
150000LC

Artix-7
360000LC

Cyclone II
68416LE
Cyclone III/LS
119088LE

Cyclone IV/E/GX
149760LE

Cyclone V
/E/GX/GS/GT
301000LE

High-end/Low-cost: X3 – X5

Technology vs. Products (Cont.)

28nm

20nm

16nm

10nm

Virtex-7

Virtex-Ultrascale Virtex-Ultrascale+

2000000LC

5541000LC

3780000LC

**Stratix-V
/E/GX/GS/GT
359200ALM**

**Stratix-10
ARM+FPU**

**Kintex-7
480000LC**

**Kintex-Ultrascale Kintex-Ultrascale+
1451000LC 1143000LC**

**Arria-IV
174000LE**

**Arria-10
ARM+FPU**

**Artix-7
360000LC**

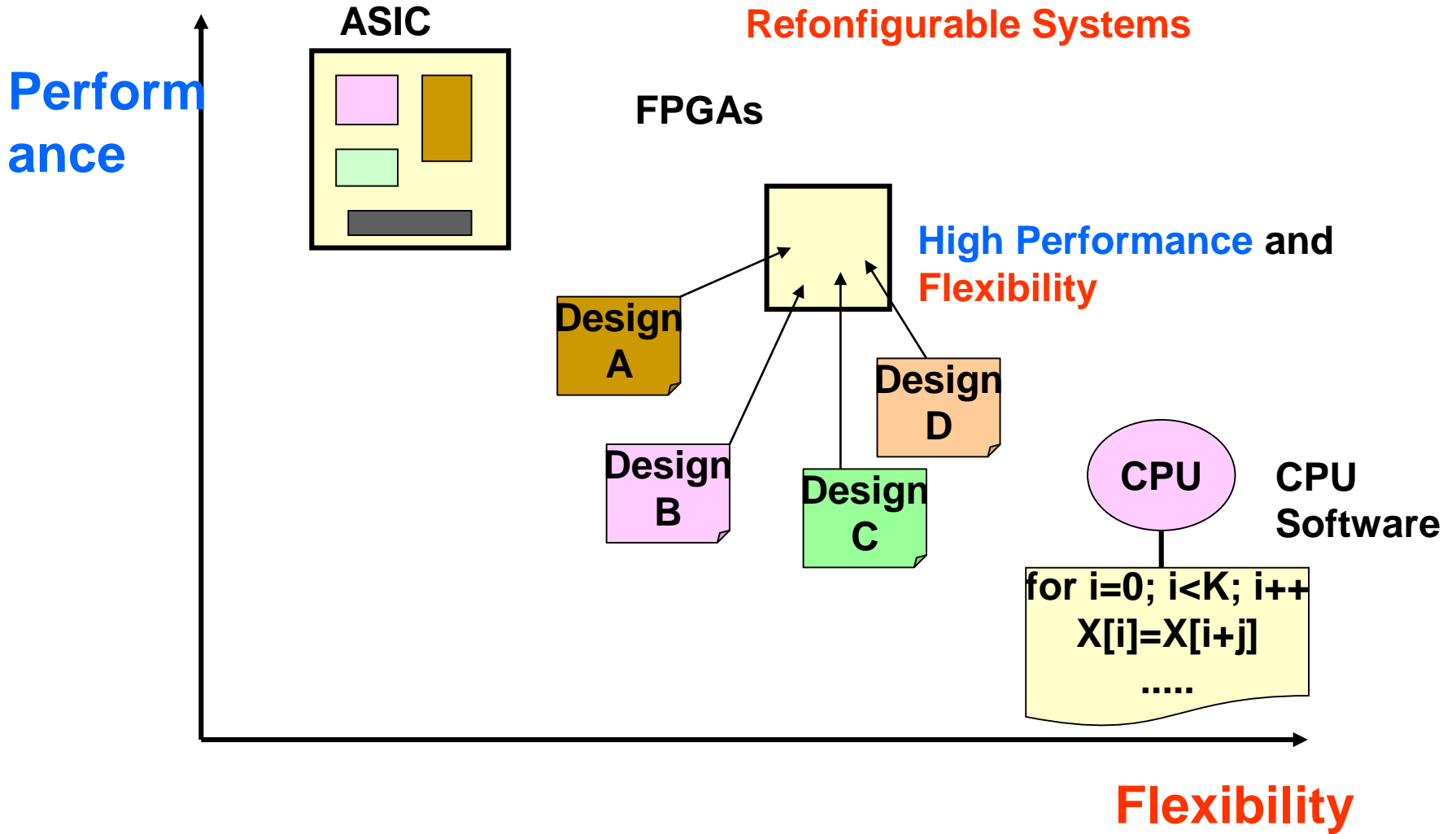
**Cyclone V
/E/GX/GS/GT
301000LE**

Design of PLDs

- Mostly designed with common HDL (Verilog-HDL, VHDL)
 - C level entry is used recently: Impulse-C, Vibado-HLS, SD-Accel(Xilinx), Open-CL, Intel-HLS(Intel)
 - Synthesis, optimization, place and route is automatically done by vendors' tools.
 - Integration and combination of tools from various vendors are used recently.
 - For large circuit, a long time is required especially for place and route.
 - Using IPs, clock/DLL adjustment is manually done.
 - Optimization techniques are different from vendors/products.
-

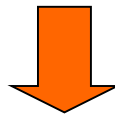
Reconfigurable System (Custom Computing Machine)

- A target algorithm is executed directly with a hardware on SRAM-style FPGA/PLDs.
 - High performance of special purpose machines.
 - High degree of flexibility of general purpose machines.
 - A completely different execution mechanism from a stored program computers.
-



How enhance the performance ?

- Performance enhancement by hardware execution itself
 - The overhead of software execution (Instruction fetch, data load to registers, and etc.)
 - The overhead of using fixed size data.
 - The overhead of using only two way branches.



However, these benefits are not so large, for embedded CPU and DSP are highly optimized.



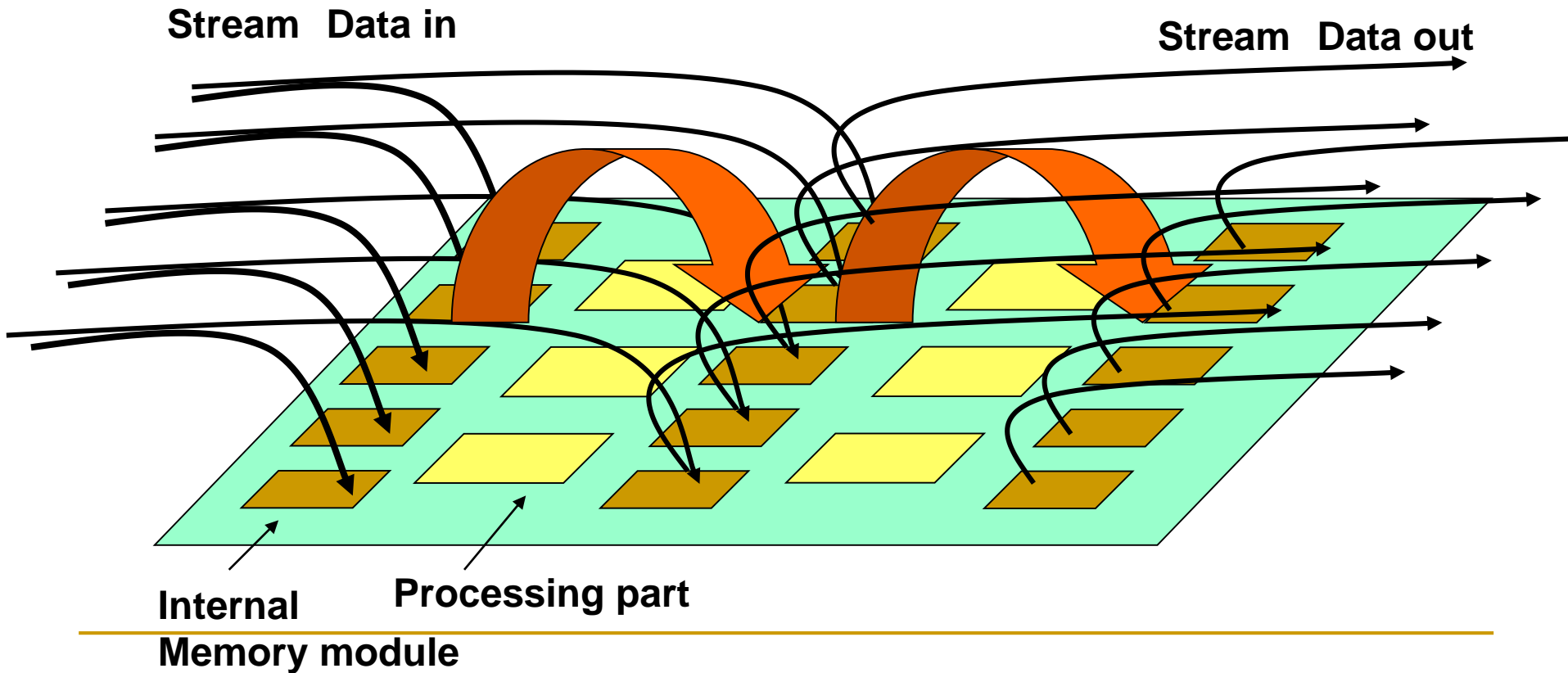
The key of performance improvement is parallel processing

Parallel processing in reconfigurable systems

- Various techniques can be used
 - SIMD execution
 - Pipelined structure
 - Systolic algorithm
 - Data driven control
 - Parallel execution other than calculation
 - Parallel data access using internal memory units
 - Parallel data transfer including I/O accesses
-

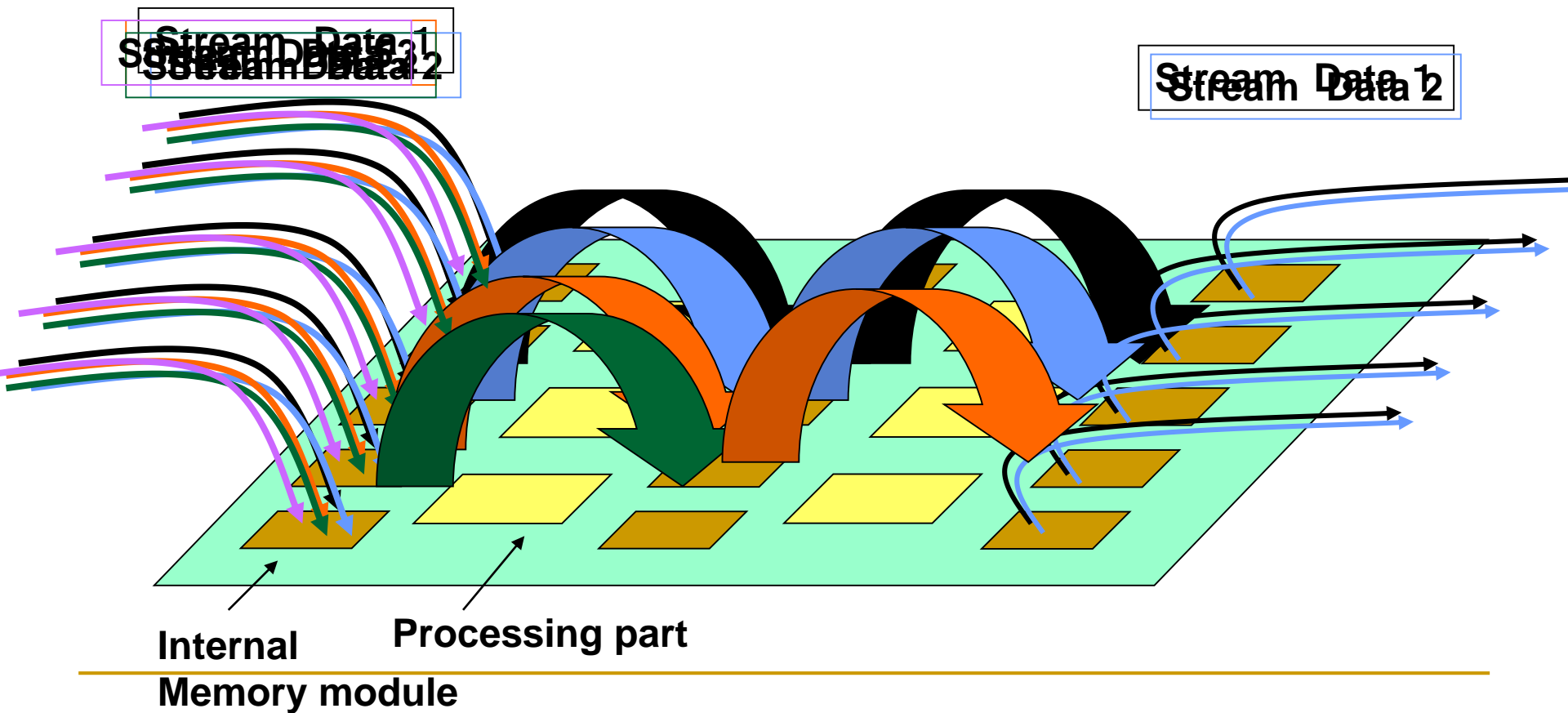
SIMD (Single Instruction-stream/ Multiple Data-stream)-like calculation

The same instruction is applied to different data stream
In Reconfigurable Systems, the operation is not required to be same
(SIMD-like calculation)

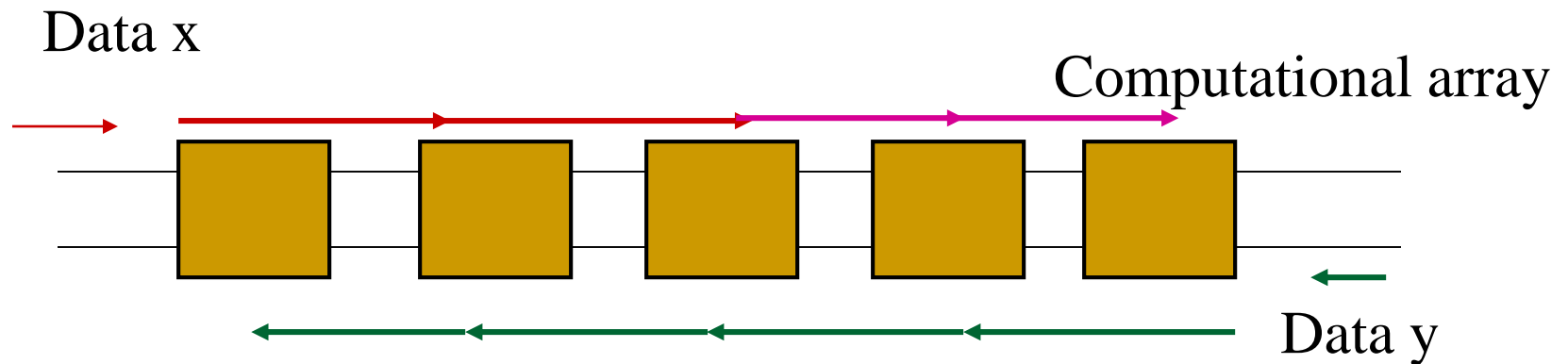


Pipelined structure

The stream is divided and inserted periodically.



Systolic Algorithm



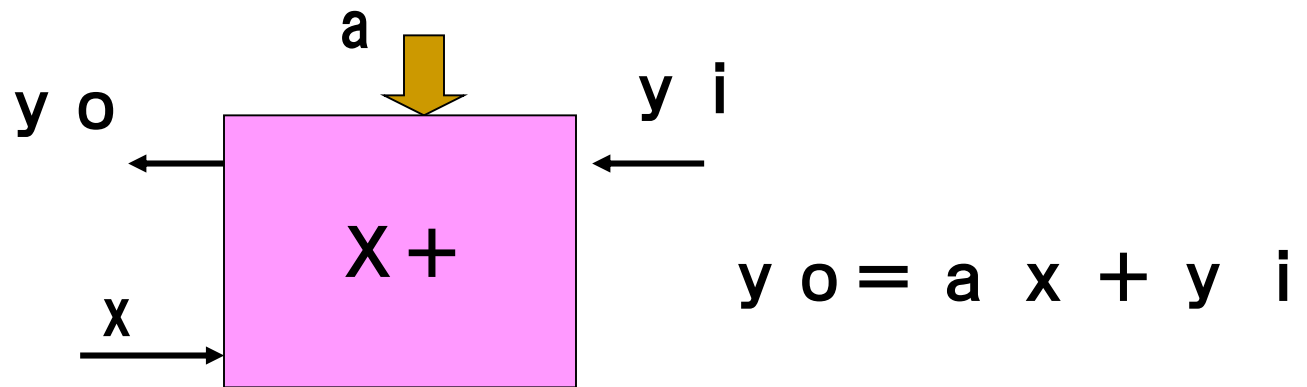
Data stream x , y are inserted with a certain interval.

When two streams meet each other, a calculation is executed.

→ Systolic: The beat of heart

Band matrix multiply $y=Ax$

$$\begin{pmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & 0 & 0 \\ a_{21} & a_{22} & a_{23} & 0 \\ 0 & a_{32} & a_{33} & a_{34} \\ 0 & 0 & a_{43} & a_{44} \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}$$



Band matrix multiply $y=Ax$



a_{23}

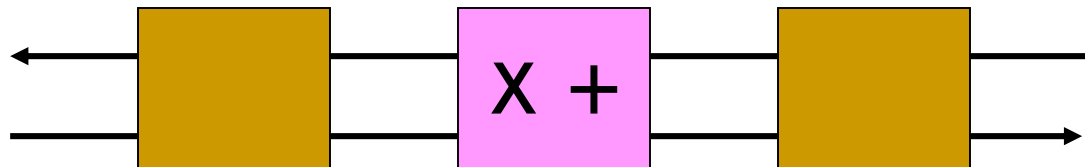
a_{32}

a_{22}

a_{12}

a_{21}

a_{11}



$$\begin{pmatrix} a_{11} & a_{12} & 0 & 0 \\ a_{21} & a_{22} & a_{23} & 0 \\ 0 & a_{32} & a_{33} & a_{34} \\ 0 & 0 & a_{43} & a_{44} \end{pmatrix}$$

x_1

Band matrix multiply $y=Ax$



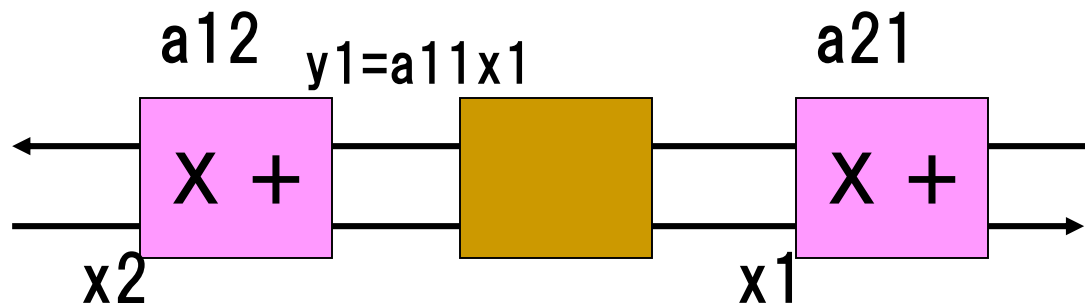
$$\begin{pmatrix} a_{11} & a_{12} & 0 & 0 \\ a_{21} & a_{22} & a_{23} & 0 \\ 0 & a_{32} & a_{33} & a_{34} \\ 0 & 0 & a_{43} & a_{44} \end{pmatrix}$$

a_{33}

a_{23}

a_{32}

a_{22}



Band matrix multiply $y=Ax$

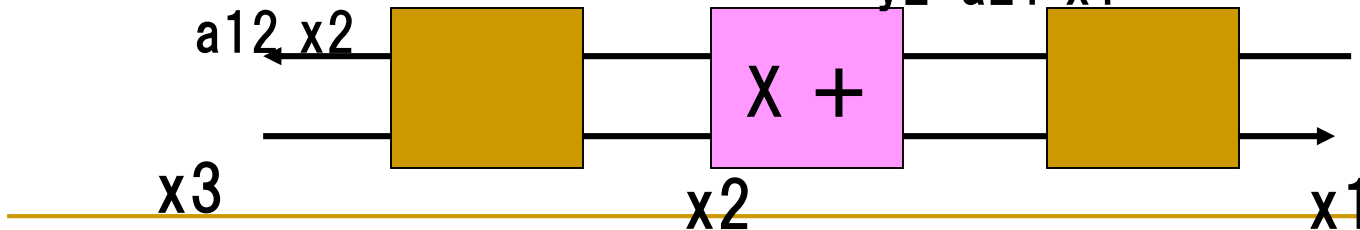


a_{34} a_{43}
 a_{23} a_{33} a_{32}
 a_{22}

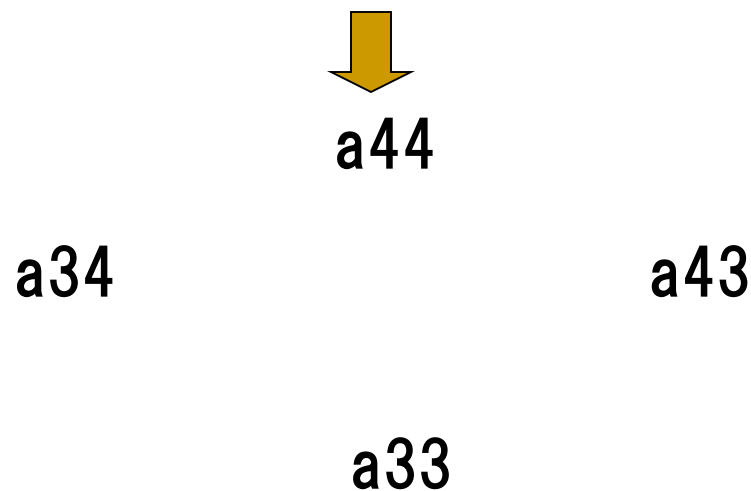
$$\begin{pmatrix} a_{11} & a_{12} & 0 & 0 \\ a_{21} & a_{22} & a_{23} & 0 \\ 0 & a_{32} & a_{33} & a_{34} \\ 0 & 0 & a_{43} & a_{44} \end{pmatrix}$$

$$y_1 = a_{11}x_1 + a_{12}x_2$$

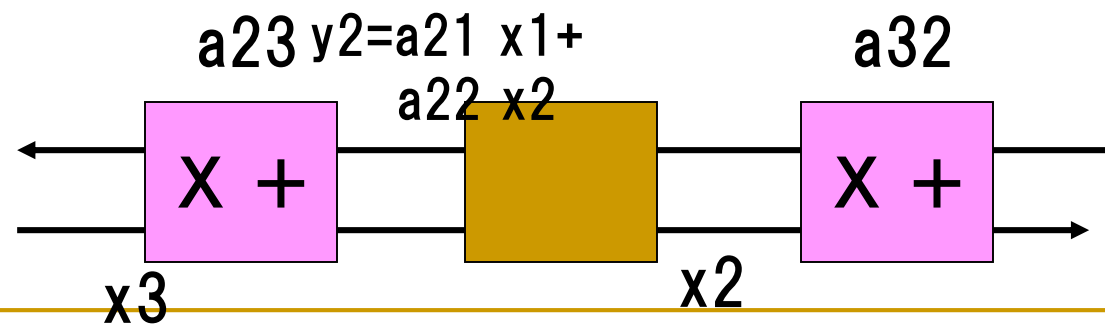
$$y_2 = a_{21}x_1$$



Band matrix multiply $y=Ax$



$$\begin{pmatrix} a_{11} & a_{12} & 0 & 0 \\ a_{21} & a_{22} & a_{23} & 0 \\ 0 & a_{32} & a_{33} & a_{34} \\ 0 & 0 & a_{43} & a_{44} \end{pmatrix}$$

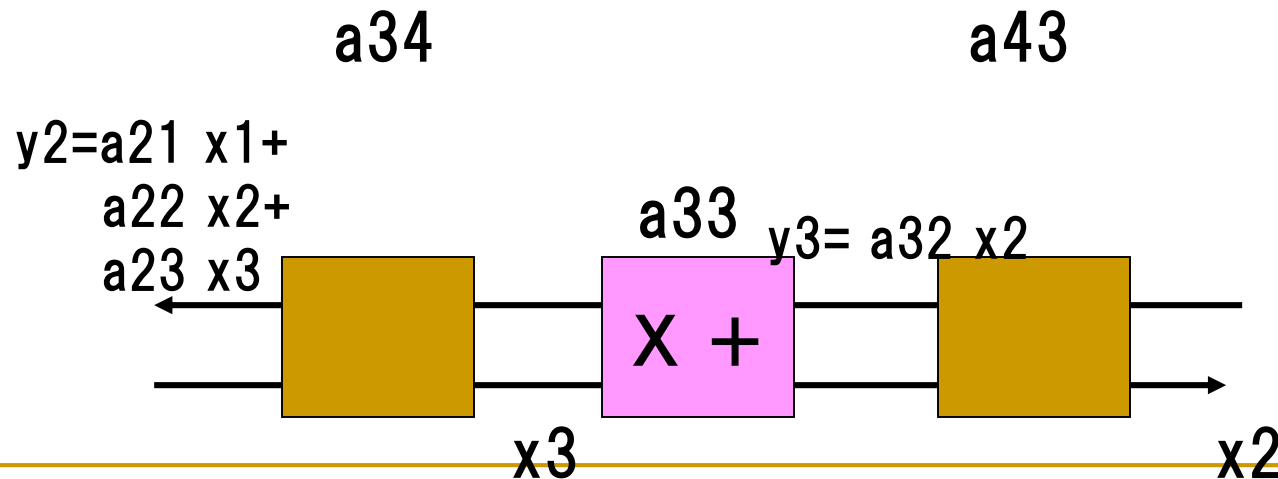


Band matrix multiply $y=Ax$

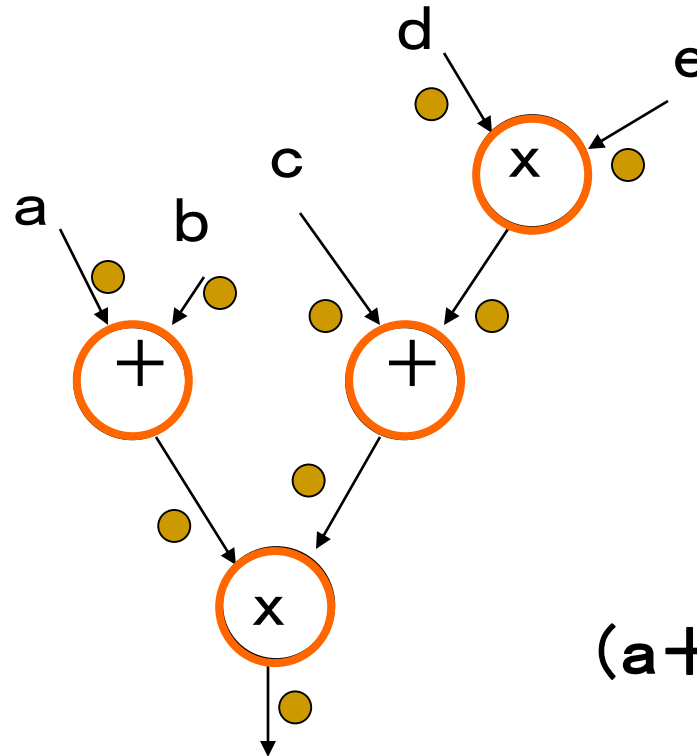


$$\begin{pmatrix} a_{11} & a_{12} & 0 & 0 \\ a_{21} & a_{22} & a_{23} & 0 \\ 0 & a_{32} & a_{33} & a_{34} \\ 0 & 0 & a_{43} & a_{44} \end{pmatrix}$$

a_{44}



Data flow algorithm

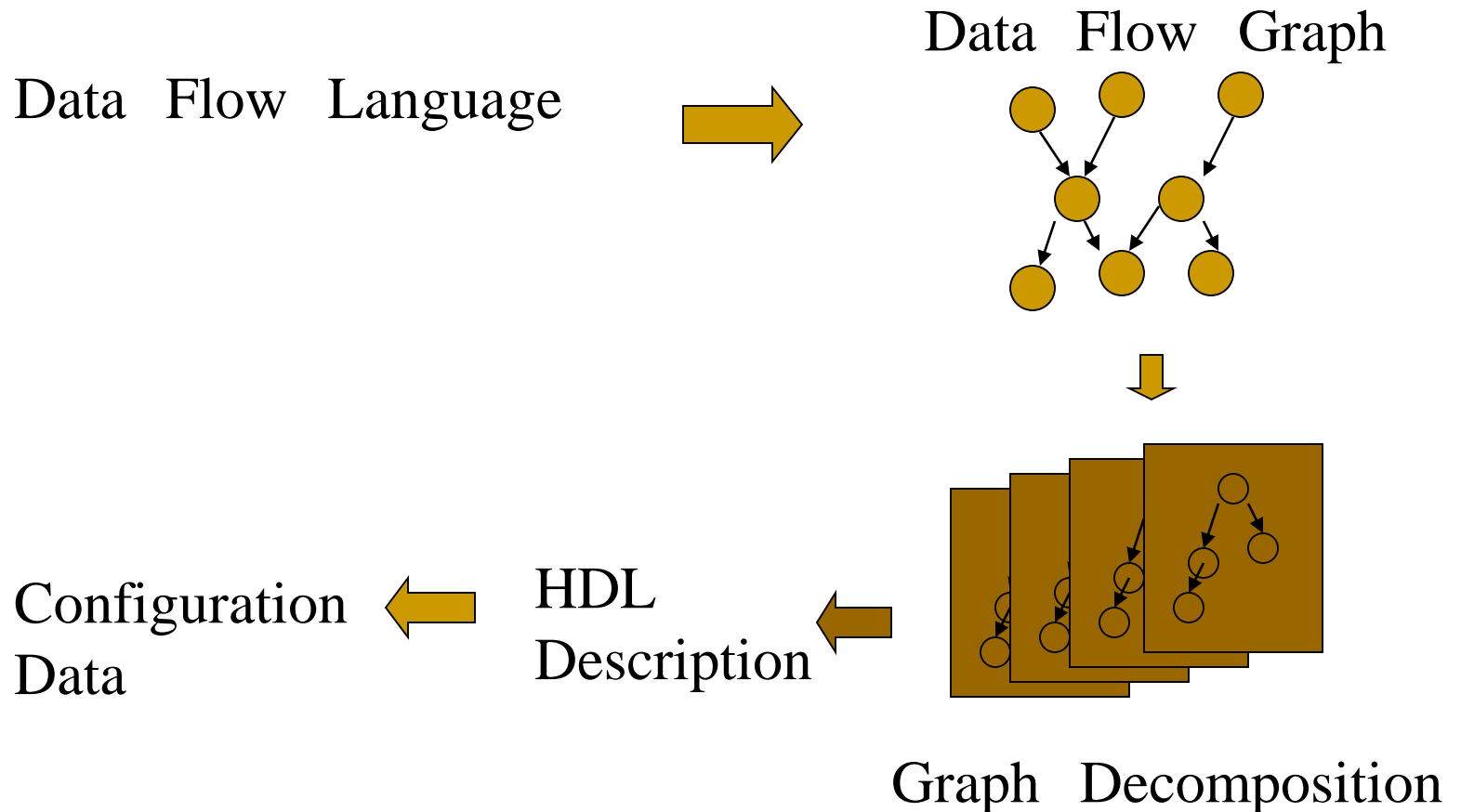


The process is activated with the available of tokens (data)

$$(a+b)x(c+(dxe))$$

The overhead of synchronization is large.

Data flow analysis and hardware generation



Suitable for automatic generation of hardware

Microsoft's Catapult

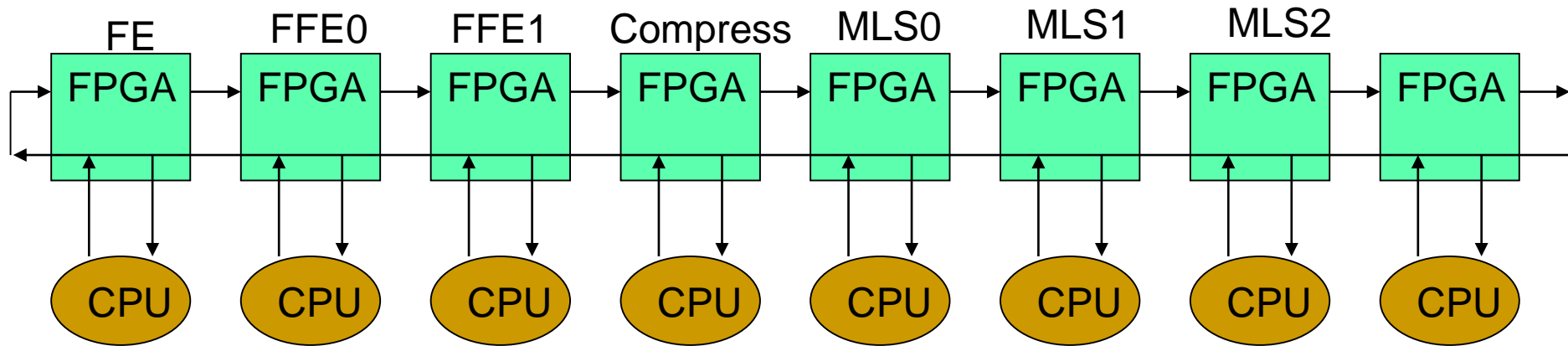
Rank computation for Web search on Bing.

Task Level Macro-Pipelining (MISD)

FE: Feature Extraction

FFE: Free Form Expression: Synthesis of feature values

MLS: Machine Learning Scoring



FPGA: Altera's Stratix V

2-Dimensional Mesh is formed (8x6) for 1 cluster.

Historical flow of computer systems

