



# Introduction to FPGA

And FPGA-related papers in recent systems conferences

Jiacheng Ma

---

# FPGAs are getting increasingly popular.

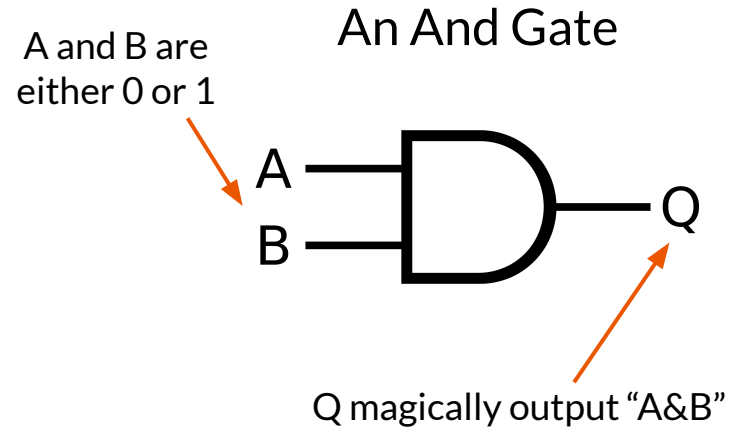
- Accelerate neural networks, graph processing, image processing, etc
- Deployed by Azure, AWS, Alibaba, etc

---

# Programing FPGAs is designing hardware.

- Unlike talks on W1 and W2
- Fundamentally different than software
  - No instruction
  - No program counter
  - No thread
  - No ...

# Let's review some elementary school physics

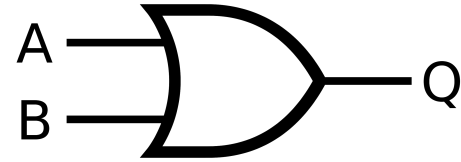


The Truth Table

A	B	Q
1	0	0
1	1	1
0	1	0
0	0	0

# Let's review some elementary school physics

An Or Gate

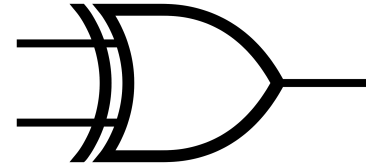


The Truth Table

A	B	Q
1	0	1
1	1	1
0	1	1
0	0	0

# Let's review some elementary school physics

## An Or Gate



## The Truth Table

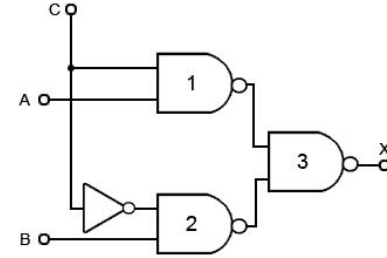
A	B	Q
1	0	1
1	1	0
0	1	1
0	0	0

The logic diagram shows a 1-bit full adder. It has three inputs: A, B, and Cin. The first 3-input OR gate has inputs A, B, and Cin, and its output is S. The first 2-input XOR gate has inputs A and B, and its output is connected to the first input of the second 2-input XOR gate. The second input of the second 2-input XOR gate is Cin. The output of the second 2-input XOR gate is Cout. The first 2-input AND gate has inputs A and B, and its output is connected to the first input of the second 2-input AND gate. The second input of the second 2-input AND gate is Cin. The output of the second 2-input AND gate is Cout.

<b>C<sub>in</sub></b>	<b>A</b>	<b>B</b>	<b>S</b>	<b>C<sub>out</sub></b>
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

# Let's review some elementary school physics

## A Selector (MUX)

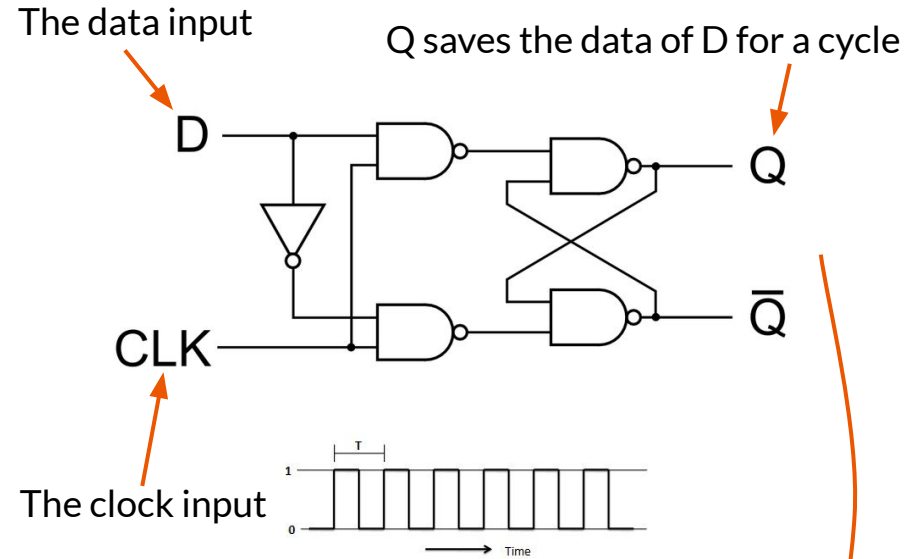


## The Truth Table

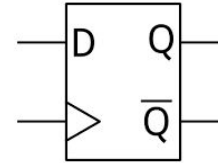
C	B	A	X
0	0	0	0
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	1



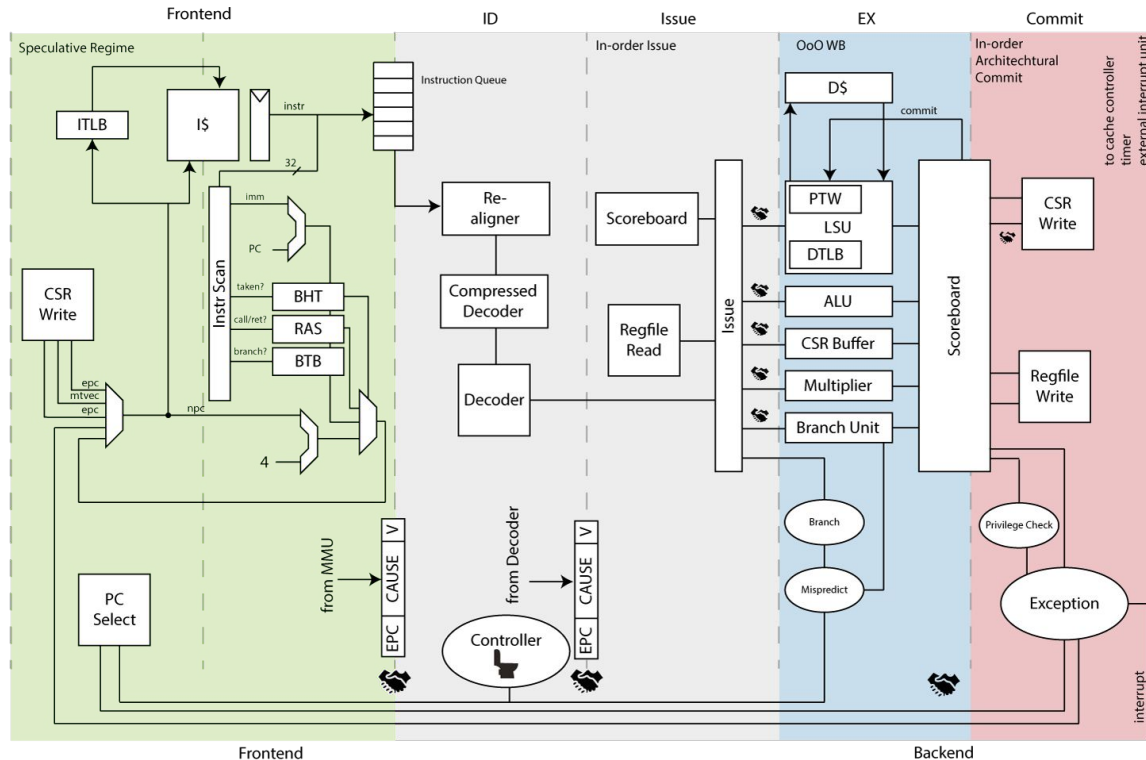
# A D-Flip-Flop



A DFF, also called a “register”



# DFFs as Pipeline Registers



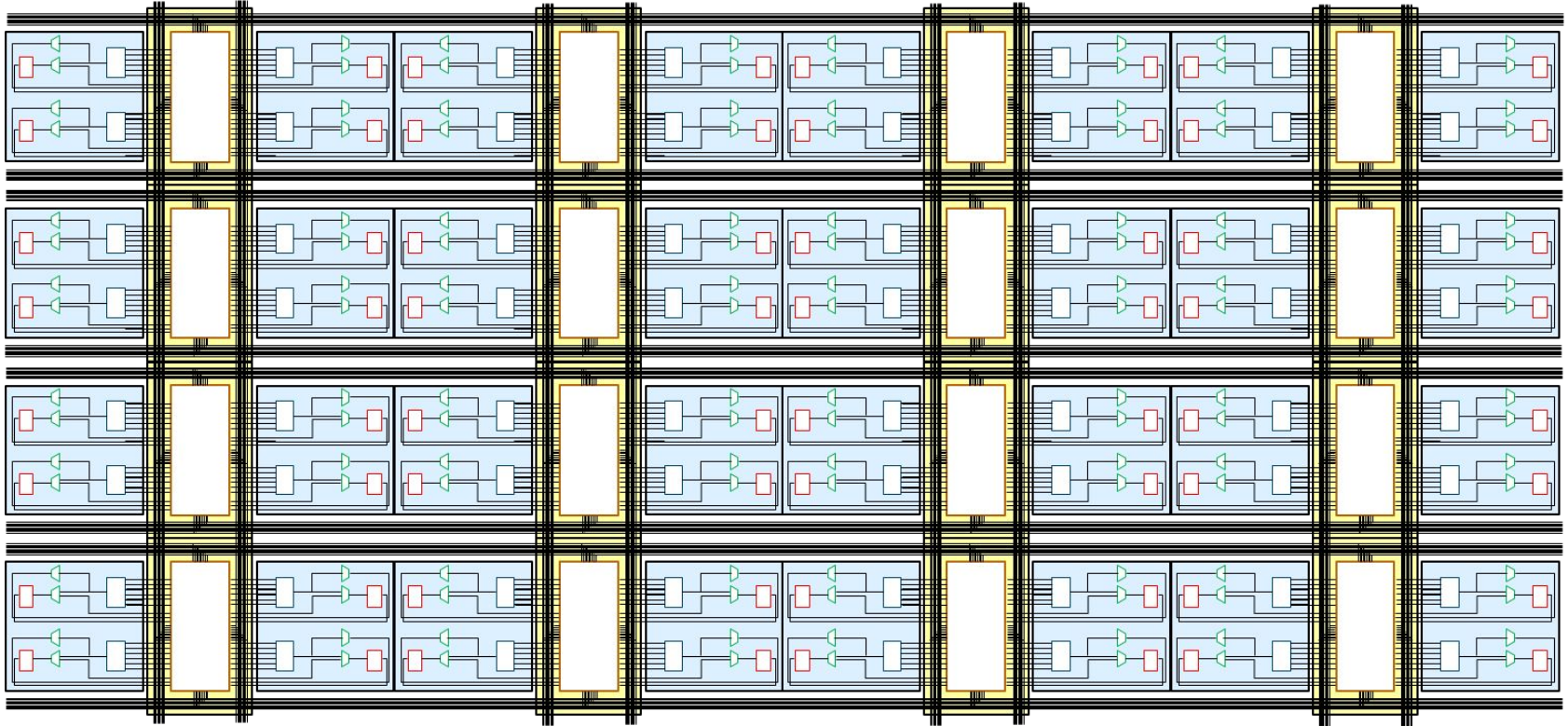
The pipeline design of  
CVA6 RISC-V processor  
<https://github.com/openhwgroup/cva6>

DFFs are used between  
different stages

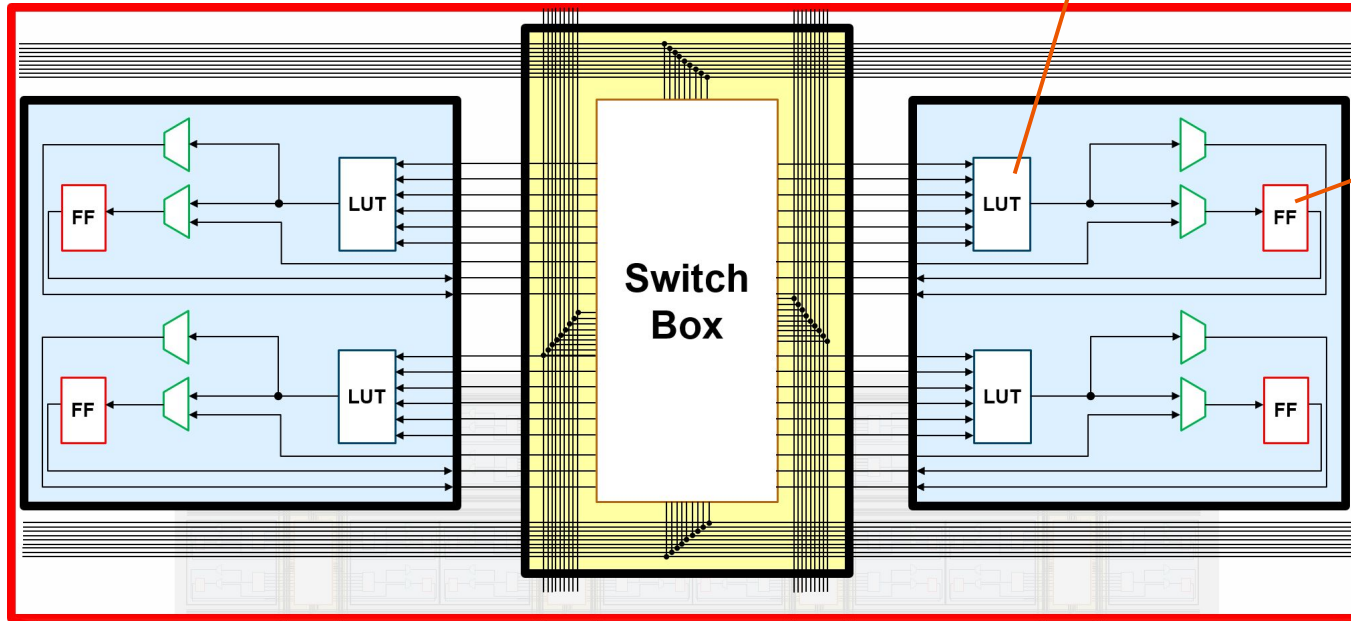
---

**But wait, where's the FPGA?**

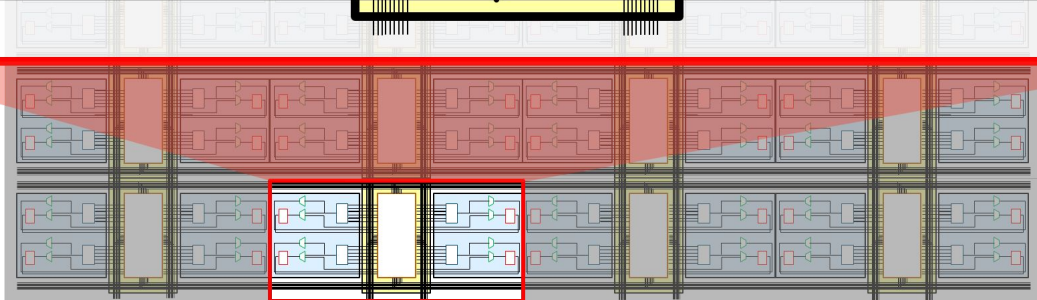
# Field Programmable Gate Array

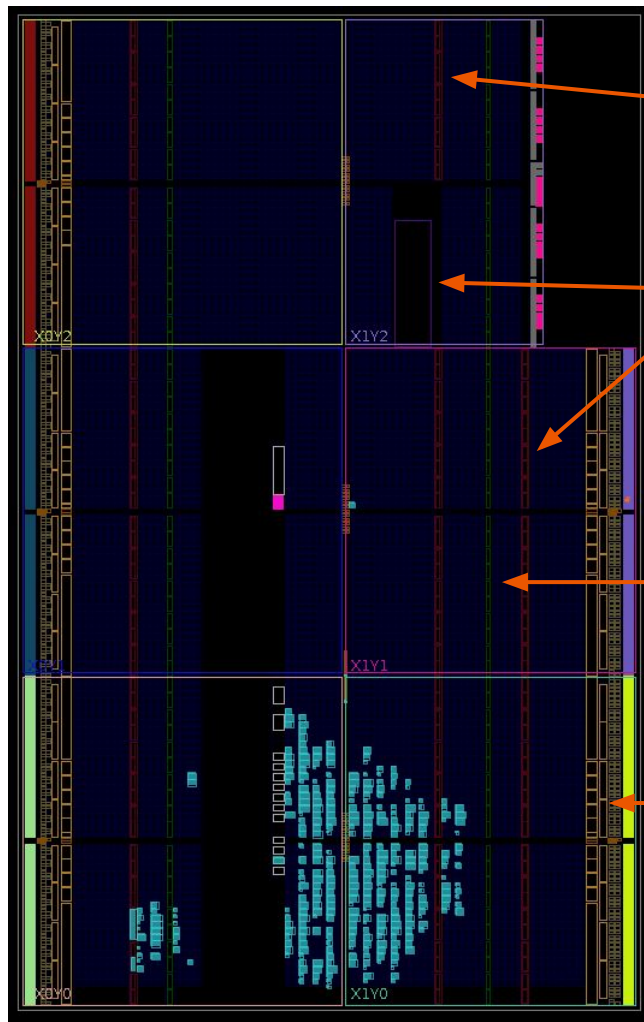


Lock Up Table: something like a truth table



Elip-Elop (register)



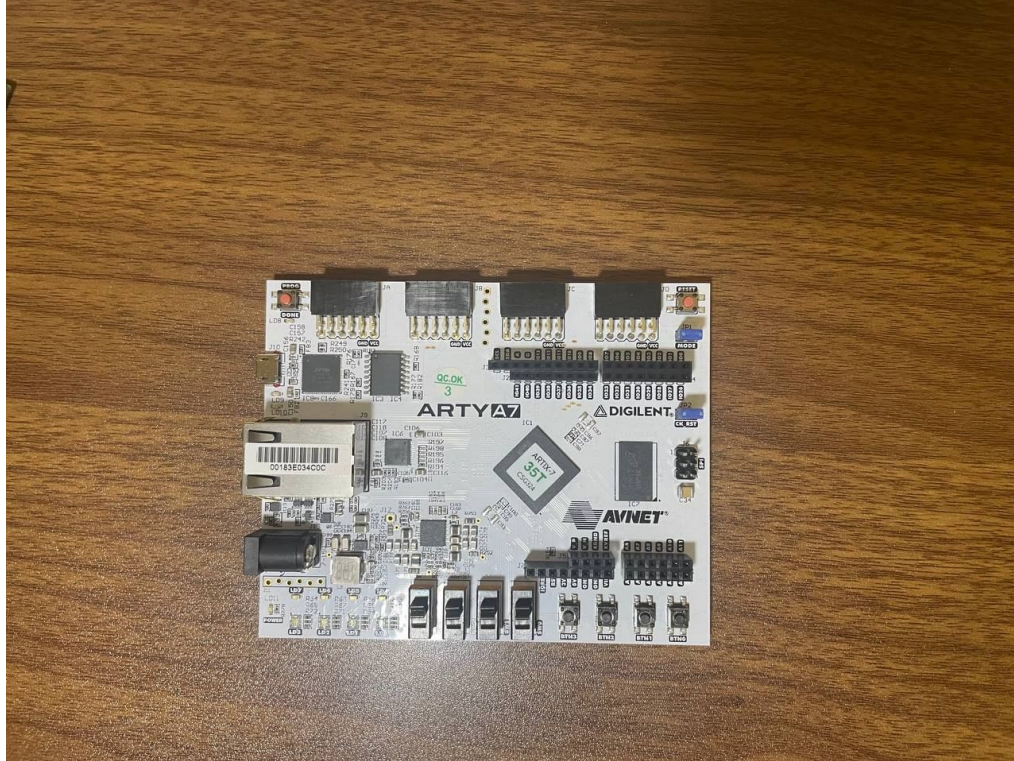


Block RAMs

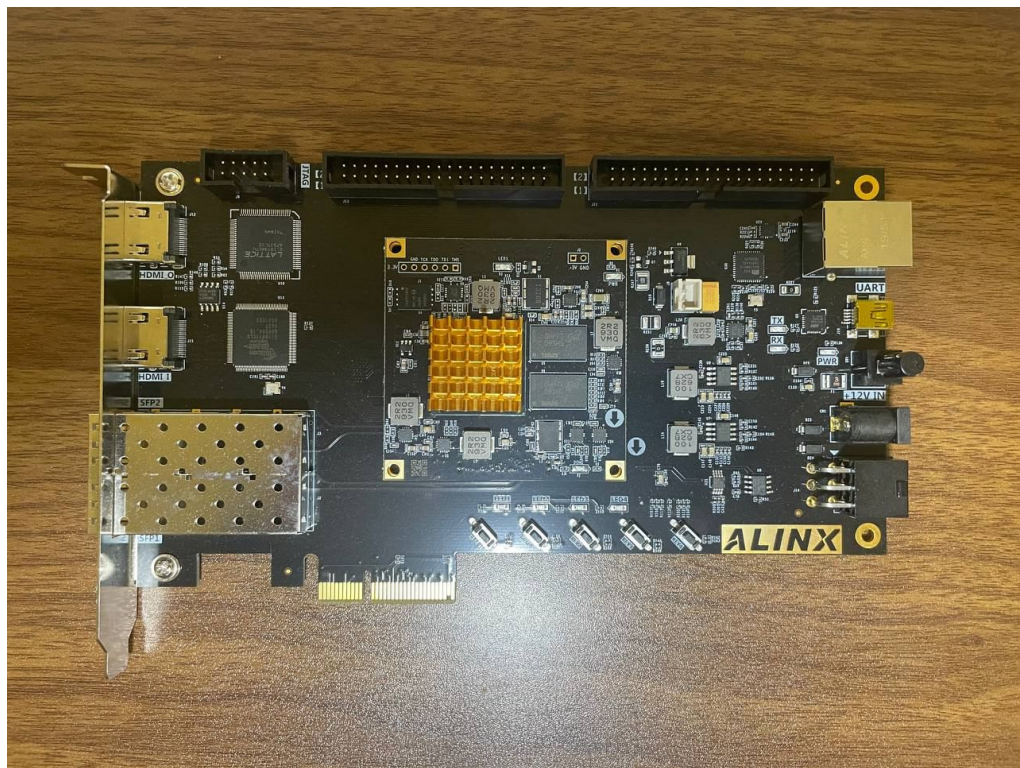
PCIe

DSPs

IO









---

**Sounds like a weird device.  
How do you program it?**

# Hardware Description Languages

Verilog, SystemVerilog, VHDL, ...

They are languages not only for FPGA, but also silicon.



# Let's use Verilog

```
1  module example(clk, valid, a, b, o1, o2);
2      input wire clk, valid;
3      input wire [1:0] a, b;
4      output reg [1:0] o1, o2;
5      reg [1:0] tmp;
6      always @(posedge clk) begin
7          if (valid) begin
8              o1 <= tmp;
9              o2 <= a + b;
10         end
11         else begin
12             o1 <= 0;
13             o2 <= 0;
14         end
15     end
16     always @(*) begin
17         tmp = a - b;
18     end
19 endmodule
```

- Module
- Always
  - With clock
  - Without clock
- Type
  - “wire”: outside always
  - “reg”: inside always
  - “logic”: everything
- Assignments
  - Blocking “=”
  - Nonblocking “<=”
    - Nonblocking → DFF
- If-statements
  - Like C
  - If → MUX
- Parallelism

# IPs



## → Close-source libraries

- ◆ Like a configurable module
- ◆ May use on-chip non-programmable resources
  - E.g., memory controller, block ram, PCIe

# Synthesis



→ A “synthesizer”

- ◆ Like a software compiler
- ◆ Verilog → a bunch of gates and wires
  - Mapping to the LUTs, FFs, and MUXs
  - A “bitstream” which configures the FPGA

# Reconfiguration



- FPGA is reconfigurable
  - ◆ Program with different bitstream
  - ◆ Only one bitstream at a time
  - ◆ No context save/restore
    - Kill and relaunch?

---

**I'm a software guy and I'm lazy.  
Verilog is too complex for me.**

# High Level Synthesis

Something that translates C/C++/OpenCL/LLVMIR to Verilog.  
So a software programmer can write hardware.

Cons: slow, resource consuming





---

**Why do we want to use FPGA?**  
**Why not use GPU?**

# Why not use GPU?

GPU works well in some tasks but not all.  
With FPGA, you can design very specialized hardware.  
Also, power consumption matters.



---

**If you need very specialized hardware,  
why not just build a silicon chip?**

# Why not build a chip?

People are not that rich.  
Also, silicon fabrication is not fast.



---

# Research on FPGA?

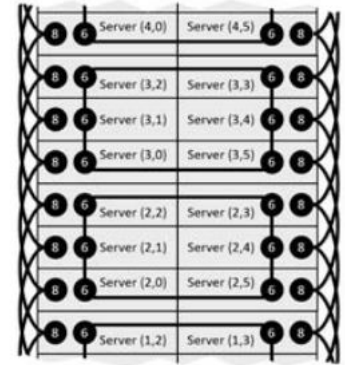
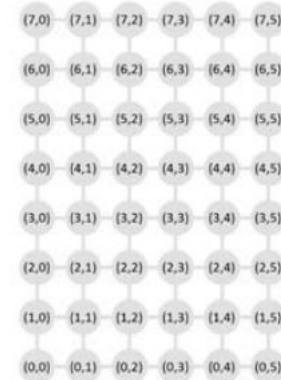
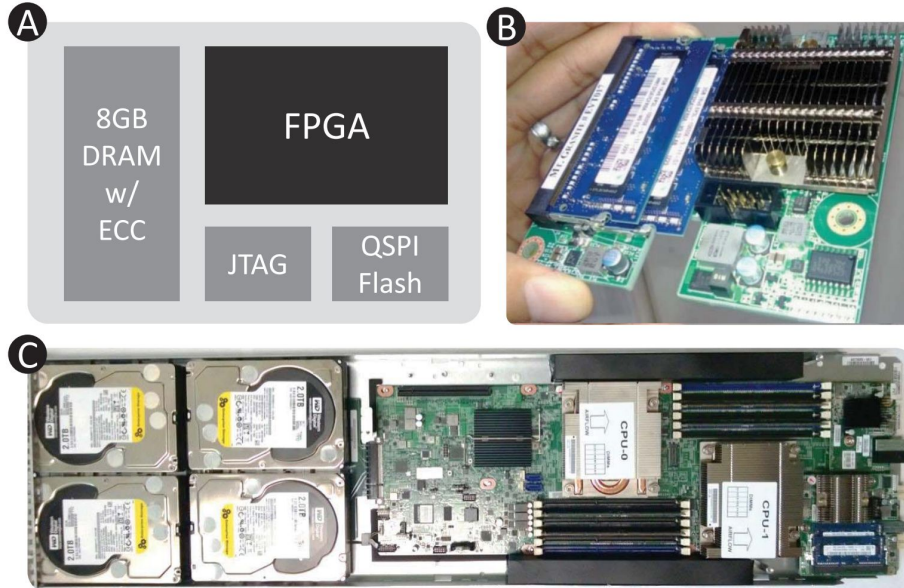


# Catapult

How does Microsoft use FPGA to make money?

# The Platform

→ M\$ is rich, they build their own boards for the FPGA



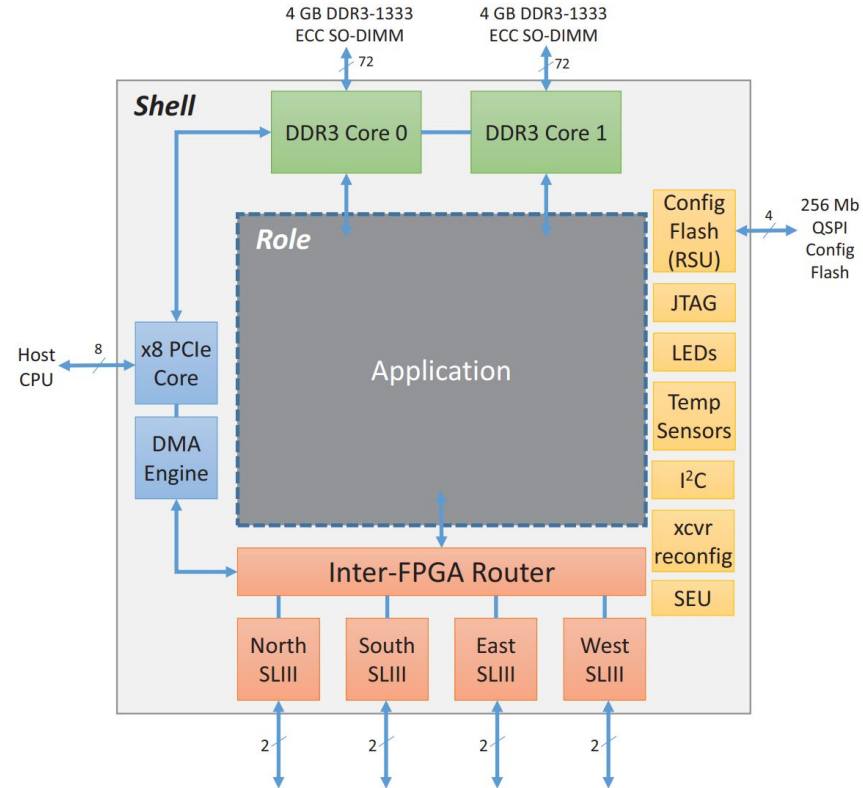
# The Shell and the Role

## → Shell

- ◆ Standard logic
- ◆ same for all applications

## → Role (~23% of resource)

- ◆ Application logic
- ◆ Partial reconfiguration





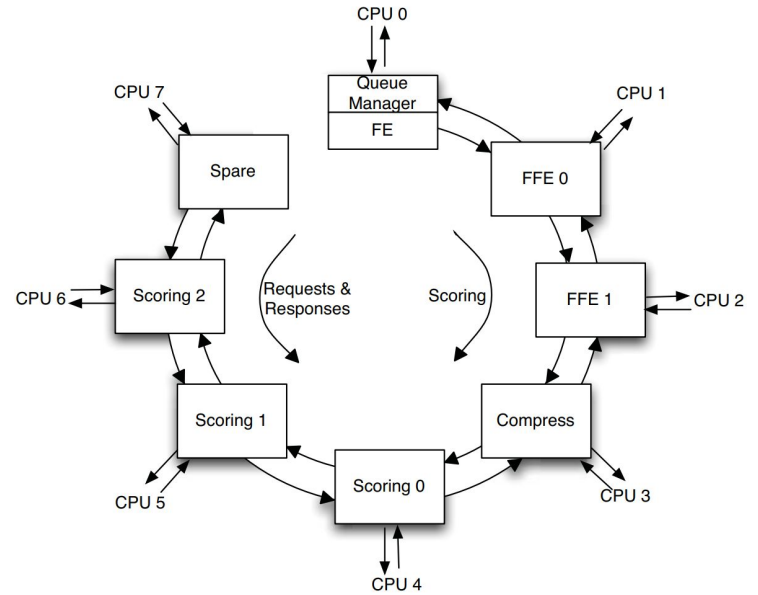
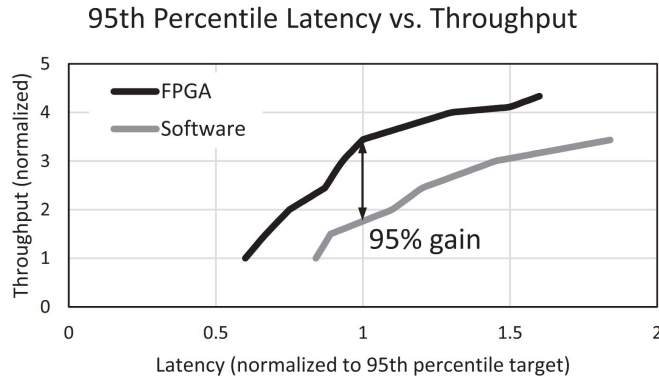
# Software Interface



- User space driver
- Per-thread buffer
  - ◆ FPGA selects the buffer to read
- Data-center software support
  - ◆ Job scheduling
  - ◆ Fault detection and recovery

# Accelerating Bing's Ranking Algorithm

- 7-FPGA pipeline for one task
  - ◆ 1 for backup, used when error occurs
- 95% speedup



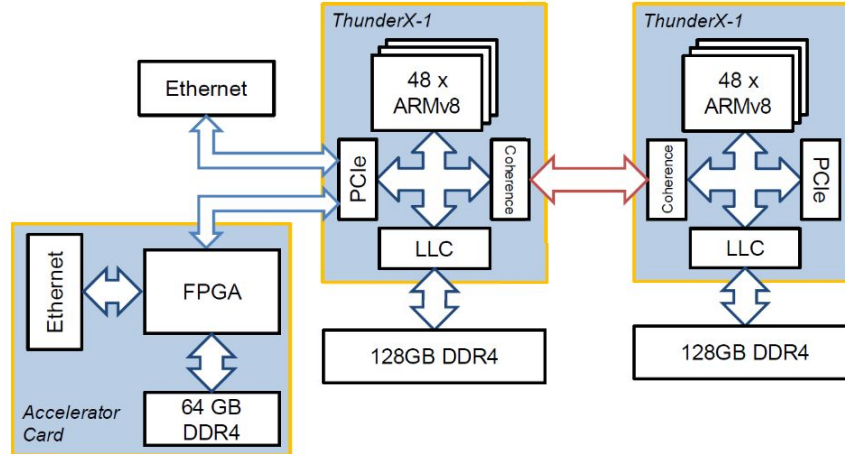


# Enzian

The future of CPU-FPGA interconnection

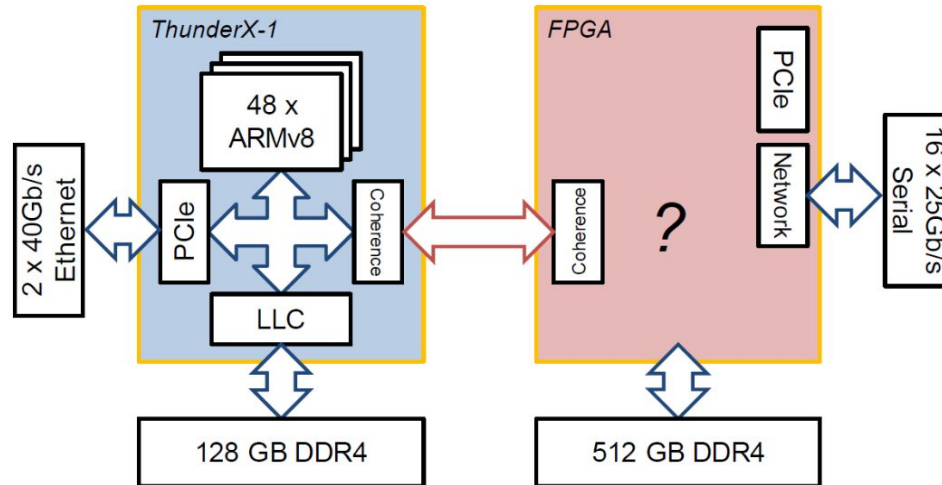
# A Traditional FPGA...

- Access CPU-side memory via PCIe
  - ◆ Relatively high latency for small transactions (while accessing memory)



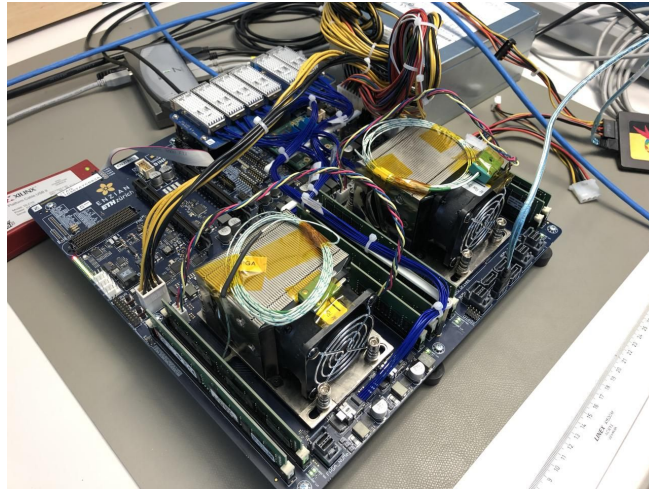
# Cache Coherent FPGA!

- Implementing cache coherent protocols on the FPGA
  - ◆ Like a NUMA node; cacheline-grained access with low latency



# Implementation of Enzian

→ 1 postdoc, a bunch of PhD students, 6 years





# AmorphOS

Now we have an OS for the FPGA...

# The Problem

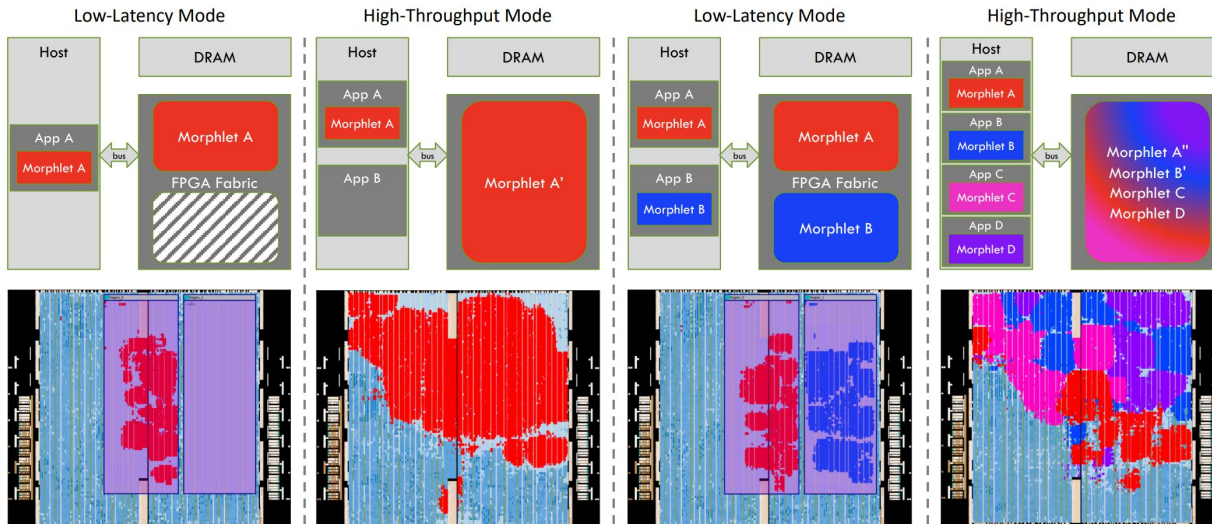


- User may not fully utilize the FPGA
  - ◆ Low resource utilization in data centers
- Different sized accelerator at different time
  - ◆ E.g., at midnight, you may want to use a small one
- Different number of instances at different time
  - ◆ E.g., Tencent's server may be more crowded during the summer break



# The solution

→ Use multiple partial-reconfigurable regions on the FPGA



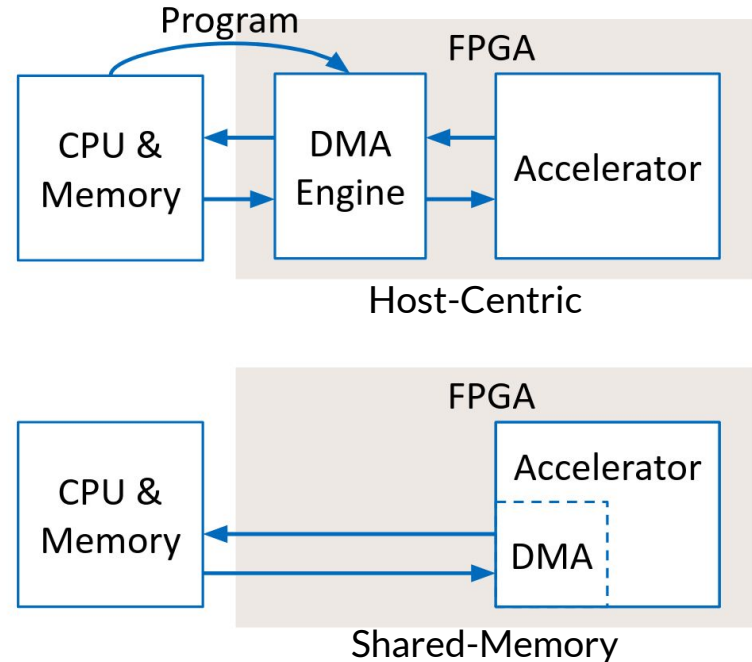


# Optimus

Exposing multiple accelerators to different VMs

# Host-Centric vs. Shared-Memory

- Whether there's a DMA engine
- Catapult/AmorphOS are for host-centric FPGAs
- How to virtualize an FPGA if it's shared-memory?



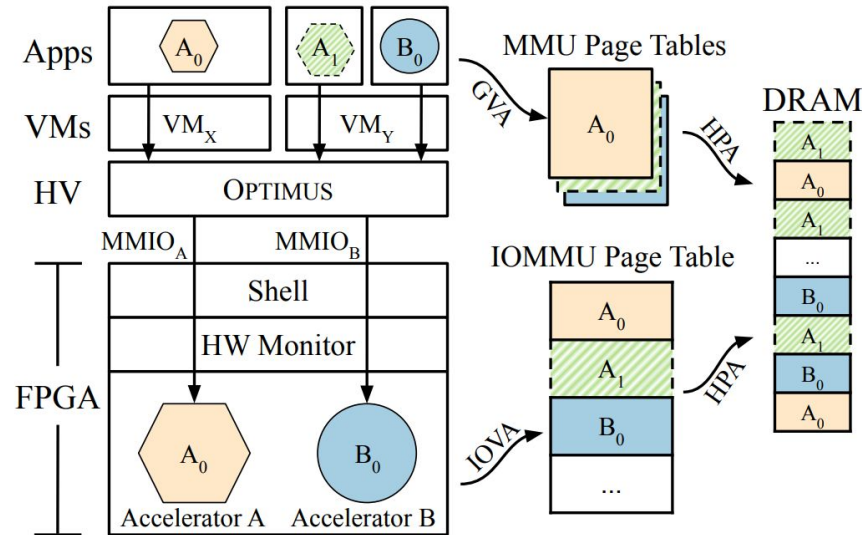
# The Problem



- No SR-IOV support
- VMs share one IOMMU page table

# The Solution

→ Divide the IOMMU page table into multiple slices



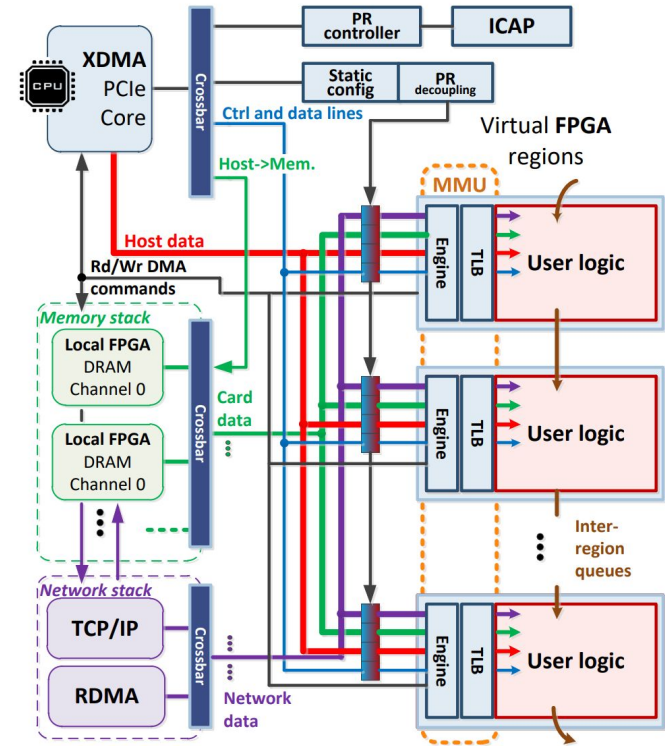


# Coyote

A more concrete implementation of FPGA OS

# The Architecture

- Per-accelerator TLB
  - ◆ Software-controlled
  - ◆ No page table
- Software-managed memory allocation



# Questions?

---