

Heterogeneous architectures and accelerators

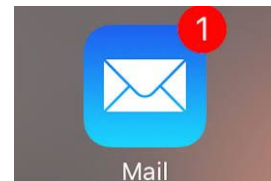
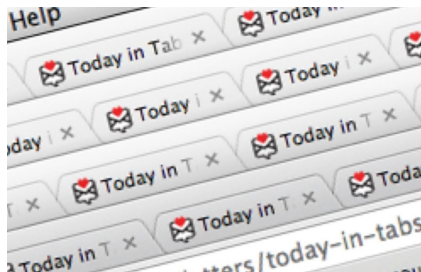


UNIVERSIDADE FEDERAL
DE MINAS GERAIS

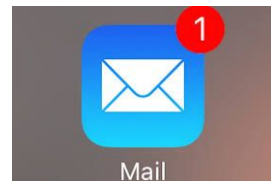
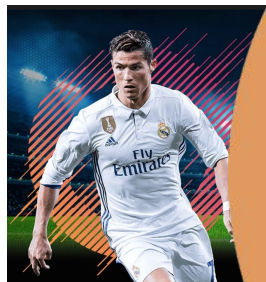
Marcelo Novaes e Rodrigo Oliveira

DCC819: Arquitetura de Computadores 2017.2

Motivation



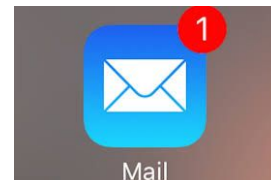
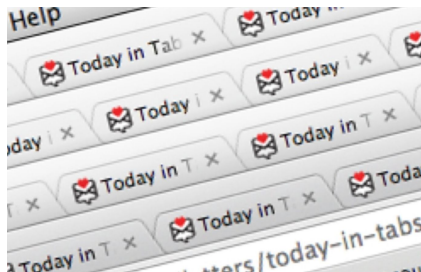
Motivation



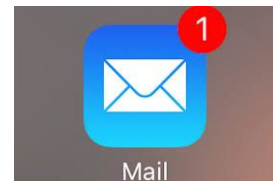
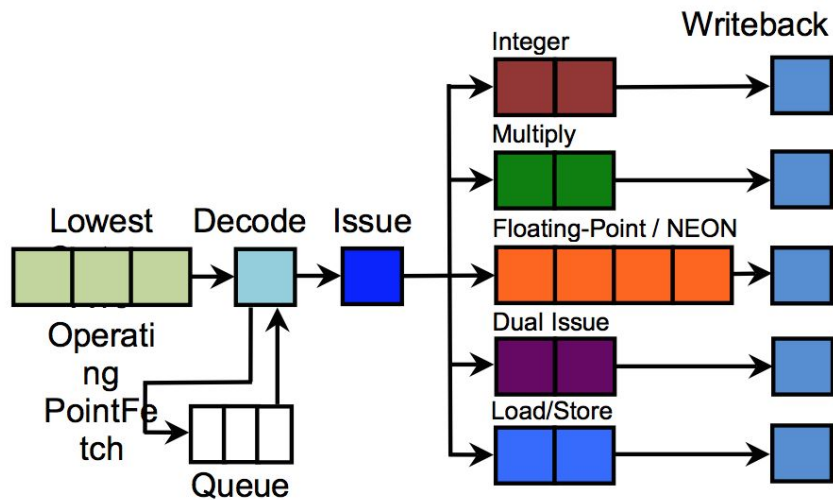
Mail

Motivation

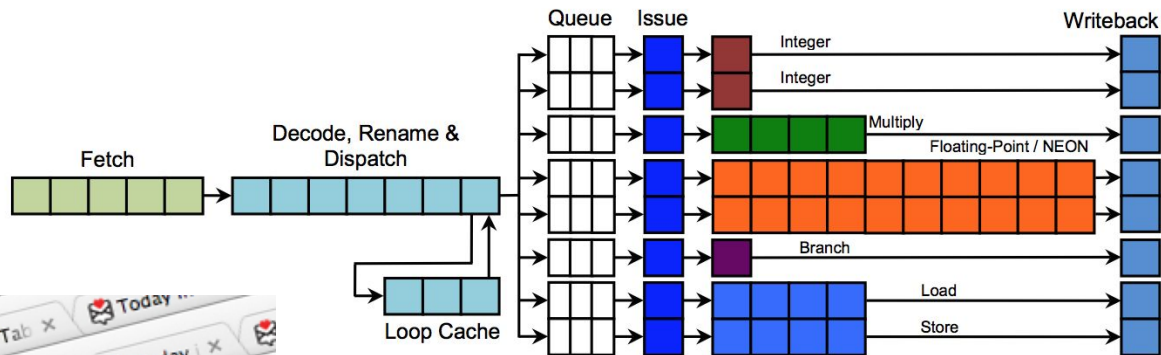
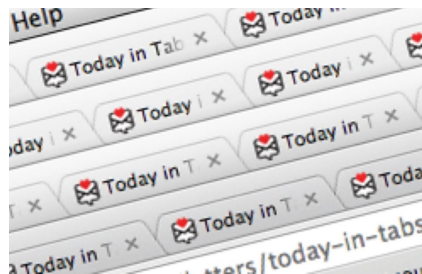
- Our use of mobile applications is heterogeneous.
 - At least two categories: High intensity tasks and Low intensity tasks.



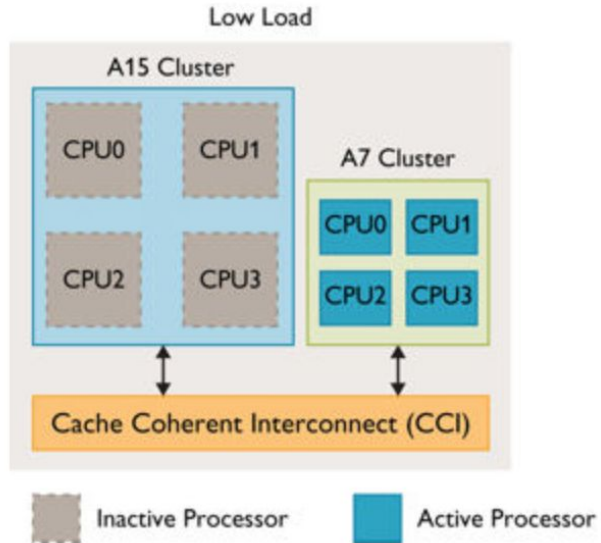
Motivation



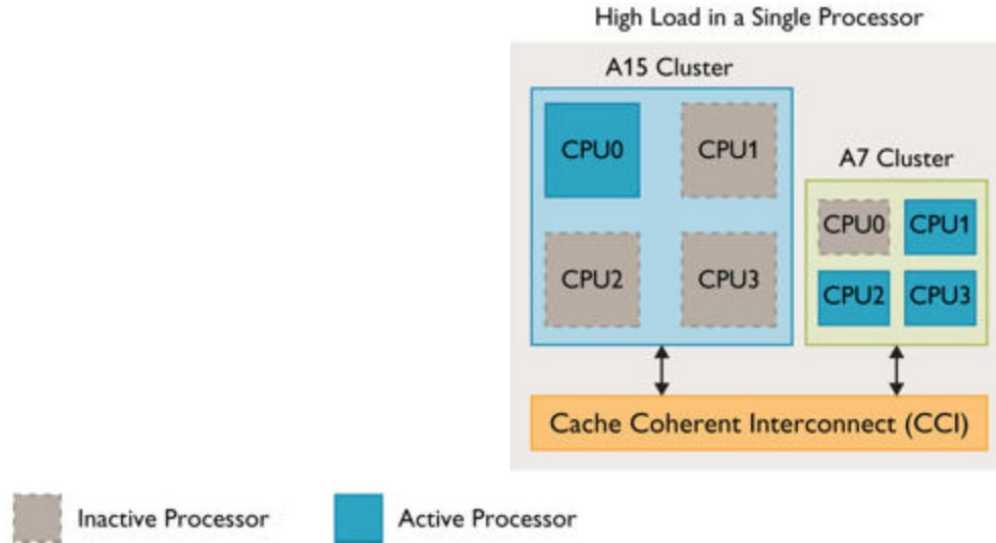
Motivation



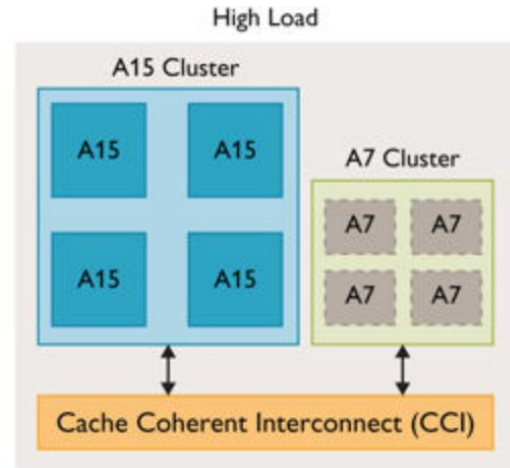
Modern Heterogeneous Multicore processors



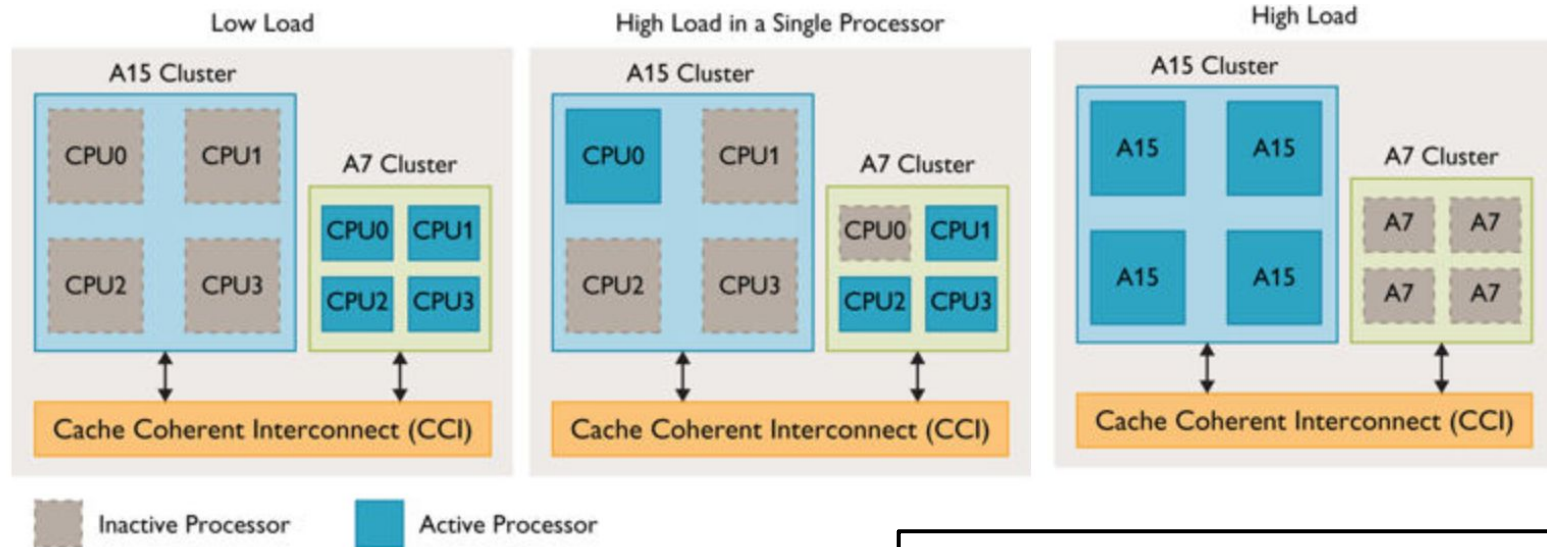
Modern Heterogeneous Multicore processors



Modern Heterogeneous Multicore processors



Modern Heterogeneous Multicore processors



Single-ISA heterogeneous multicore

Existing software can be **transparently migrated** one from another.

Big.LITTLE Architectures



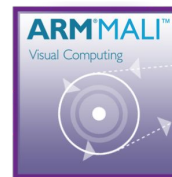
ARM MALI™
Visual Computing

Mali-300
Mali-400
Mali-450

Mali-T622
Mali-T624
Mali-T628
Mali-T678
Mali-T760

Mali-V500

ARM big.LITTLE™
Processor Technology



ARM CORELINK™ Coherent System IP
Processor System IP

ARM ARTISAN™ Physical IP

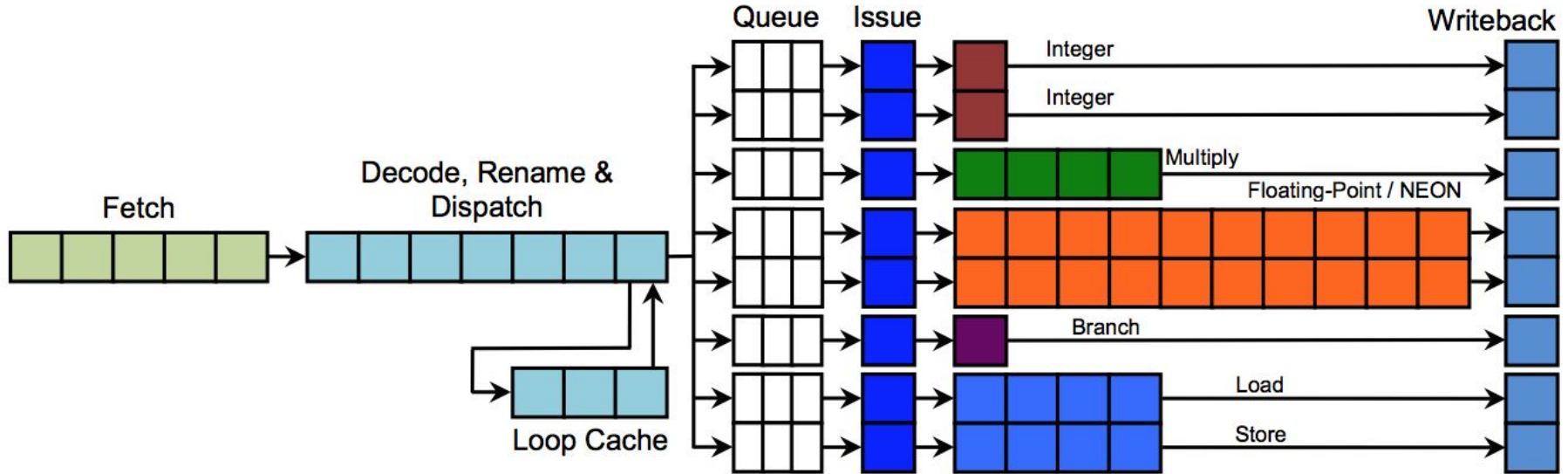
ARM CORTEX®
Processor Technology

Cortex-A57
Cortex-A53
Cortex-A15
Cortex-A9
Cortex-A7
Cortex-A5

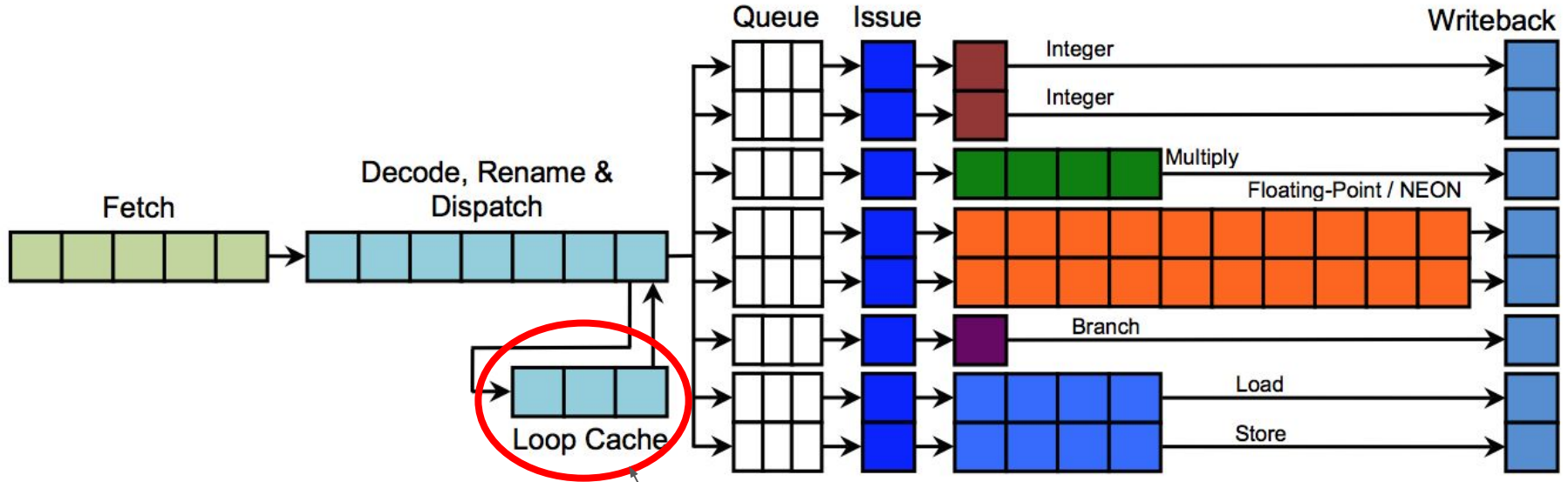
Cortex-R7
Cortex-R5
Cortex-R4

Cortex-M4
Cortex-M3
Cortex-M0

ARM Cortex-A15 (Big)

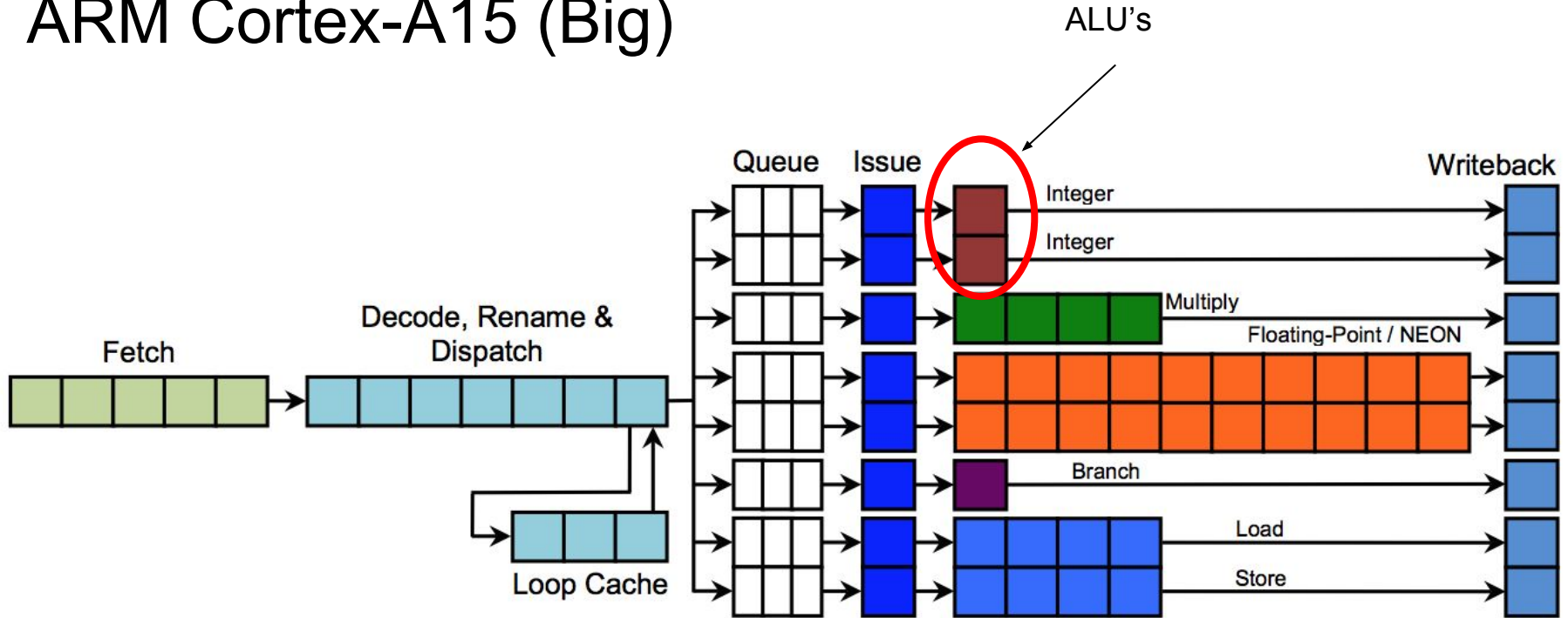


ARM Cortex-A15 (Big)

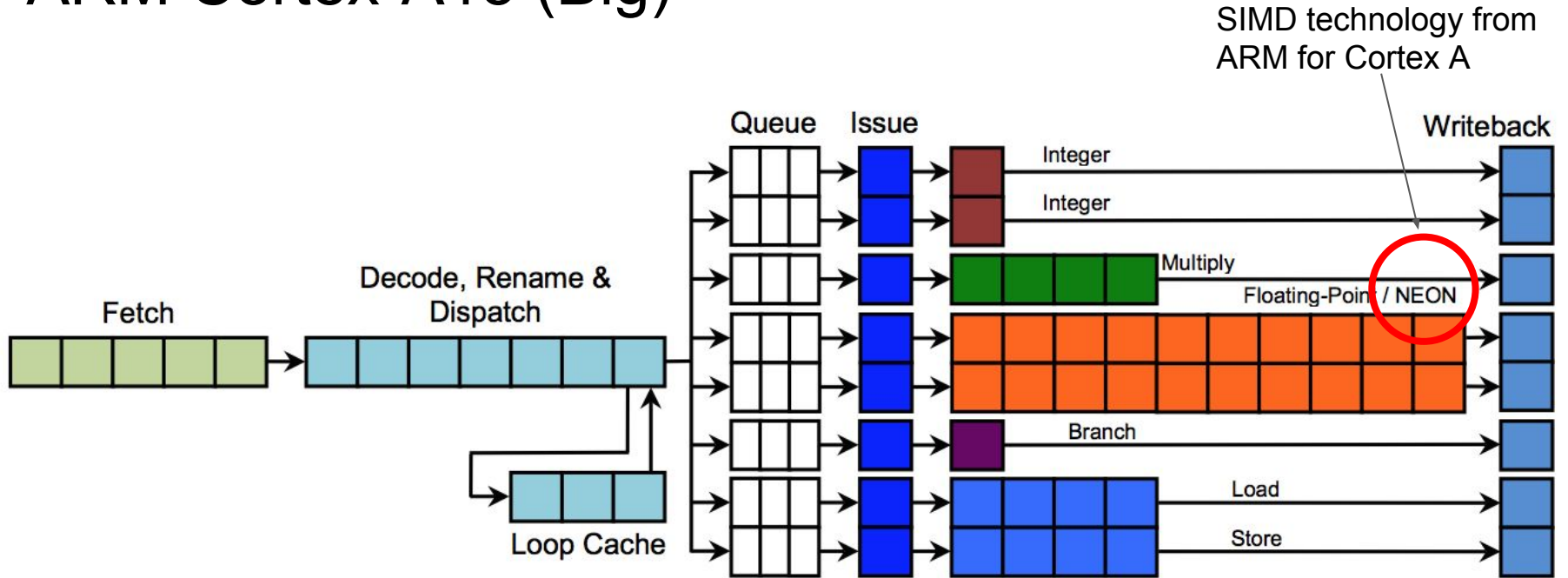


It holds instructions of small loops in order to do a Fetch faster than doing from L1 cache.

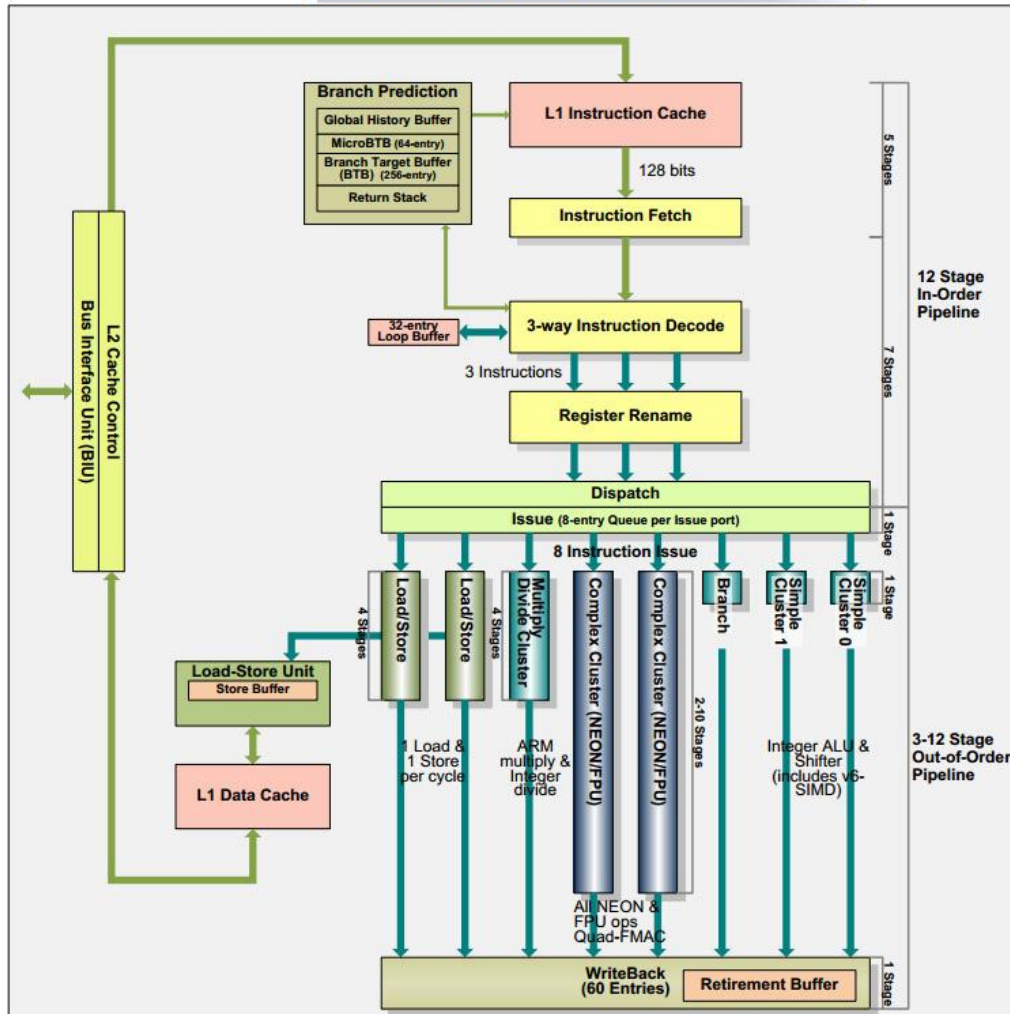
ARM Cortex-A15 (Big)



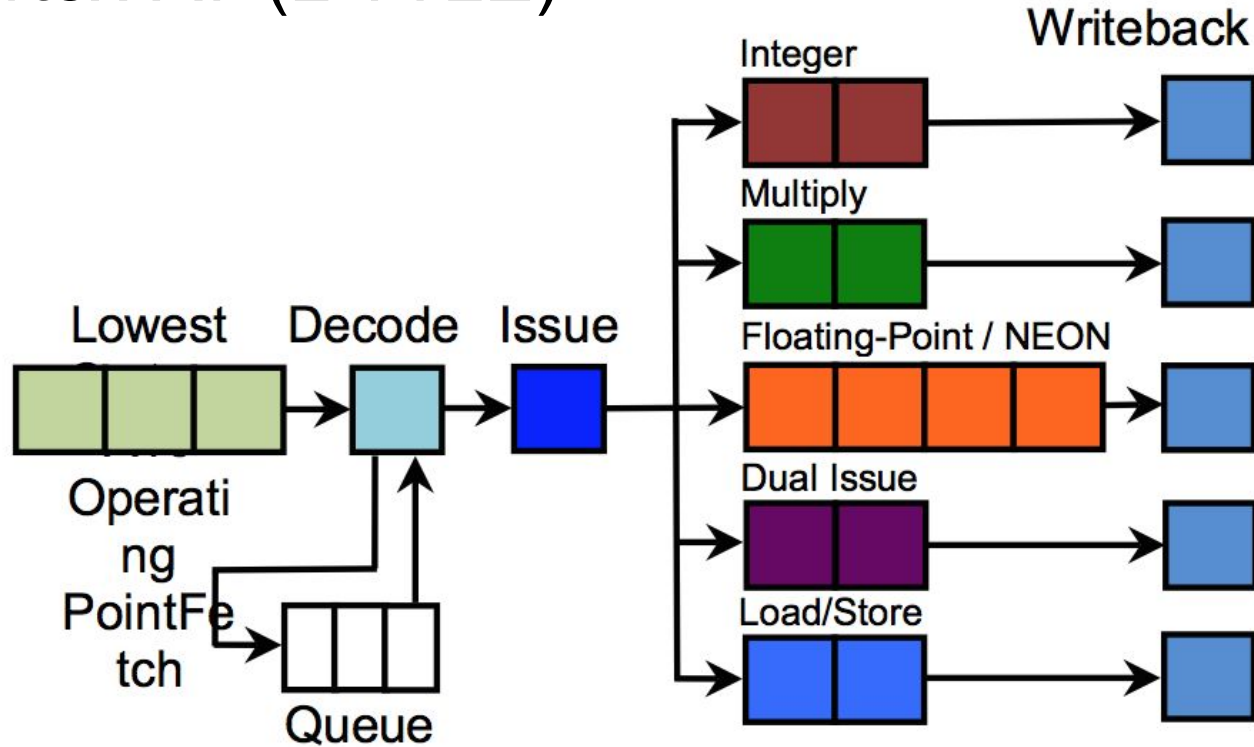
ARM Cortex-A15 (Big)



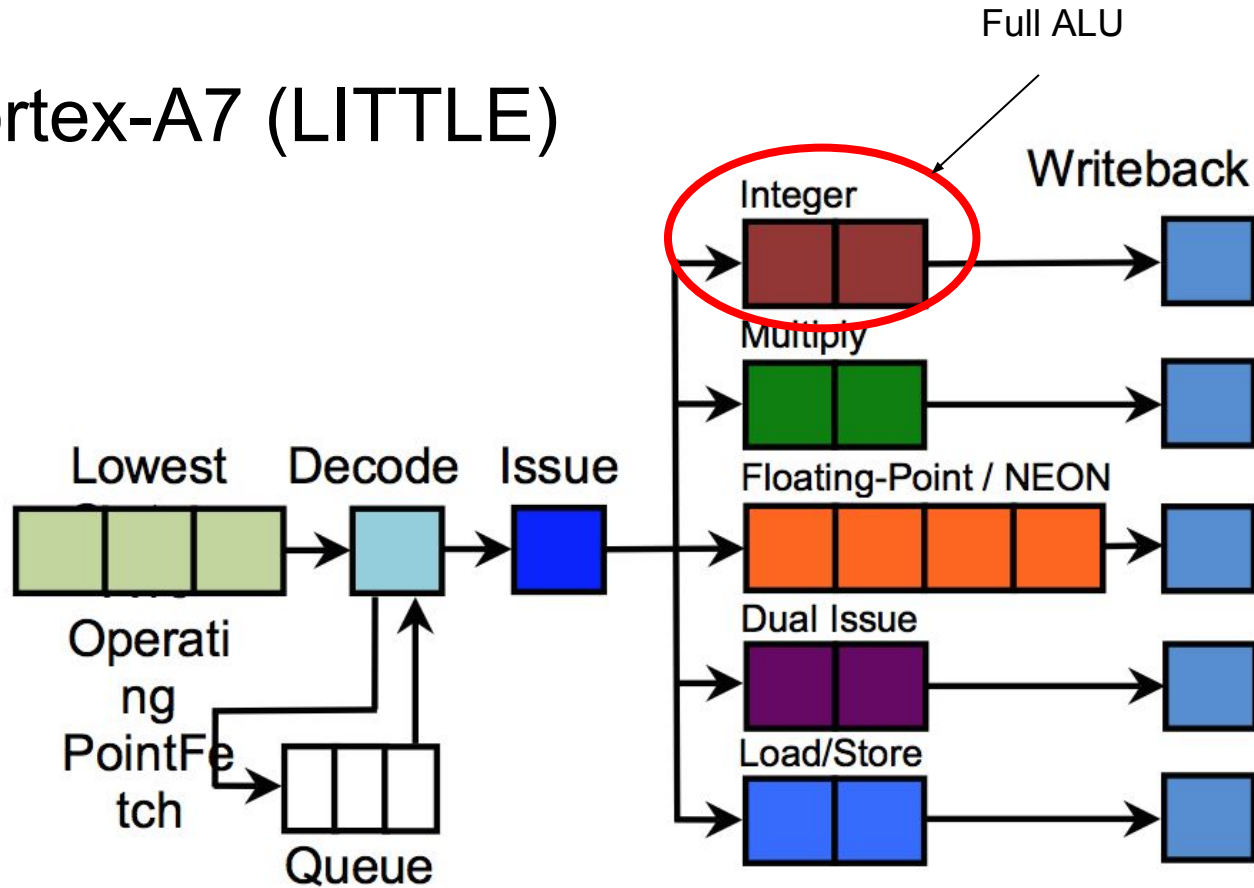
ARM Cortex-A15 Block Diagram



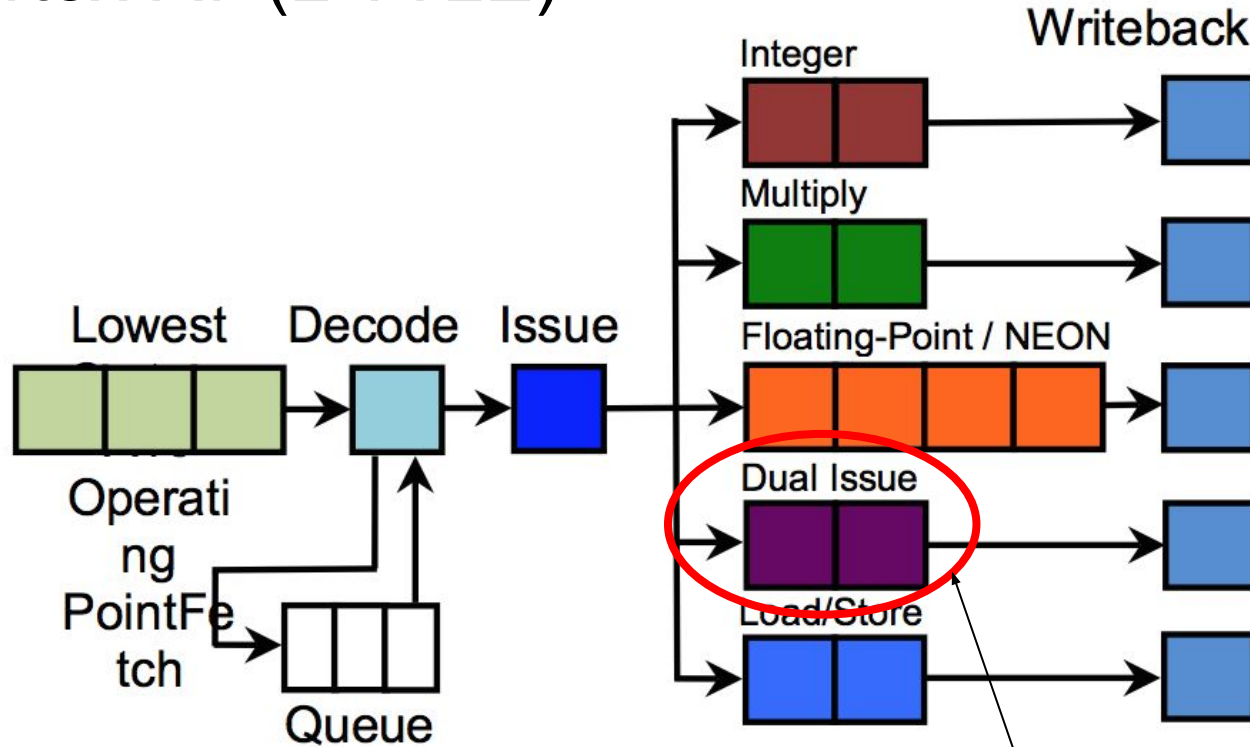
ARM Cortex-A7 (LITTLE)



ARM Cortex-A7 (LITTLE)



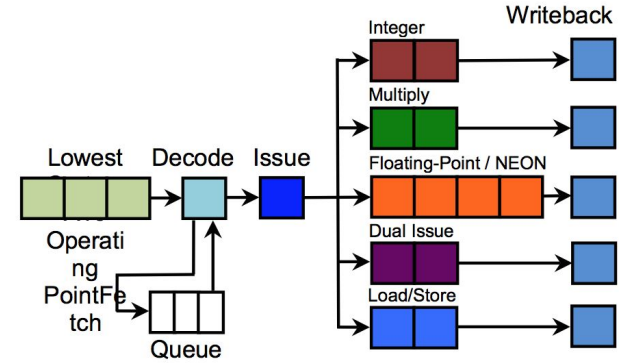
ARM Cortex-A7 (LITTLE)



"Partial" ALU (does not do all integers operations)

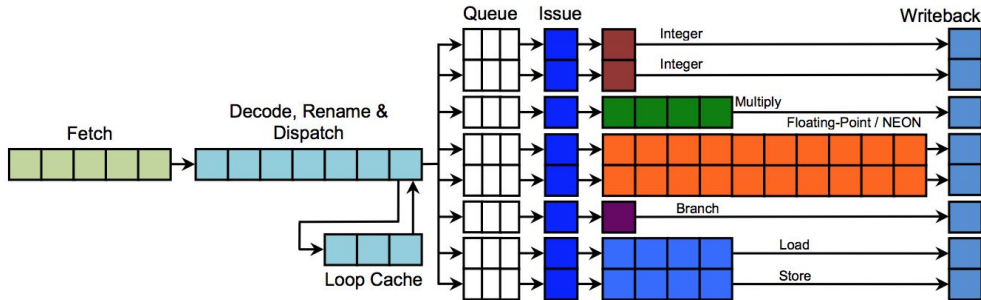
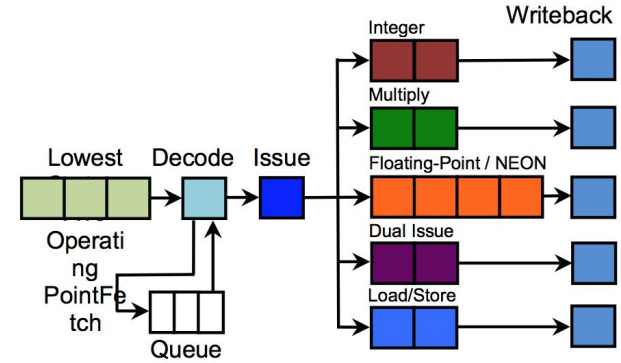
ARM Cortex-A7 (LITTLE) vs Cortex-A15 (Big)

- In-order
- Non-symmetric Dual Issue
- Processor pipeline between 8-10 stages



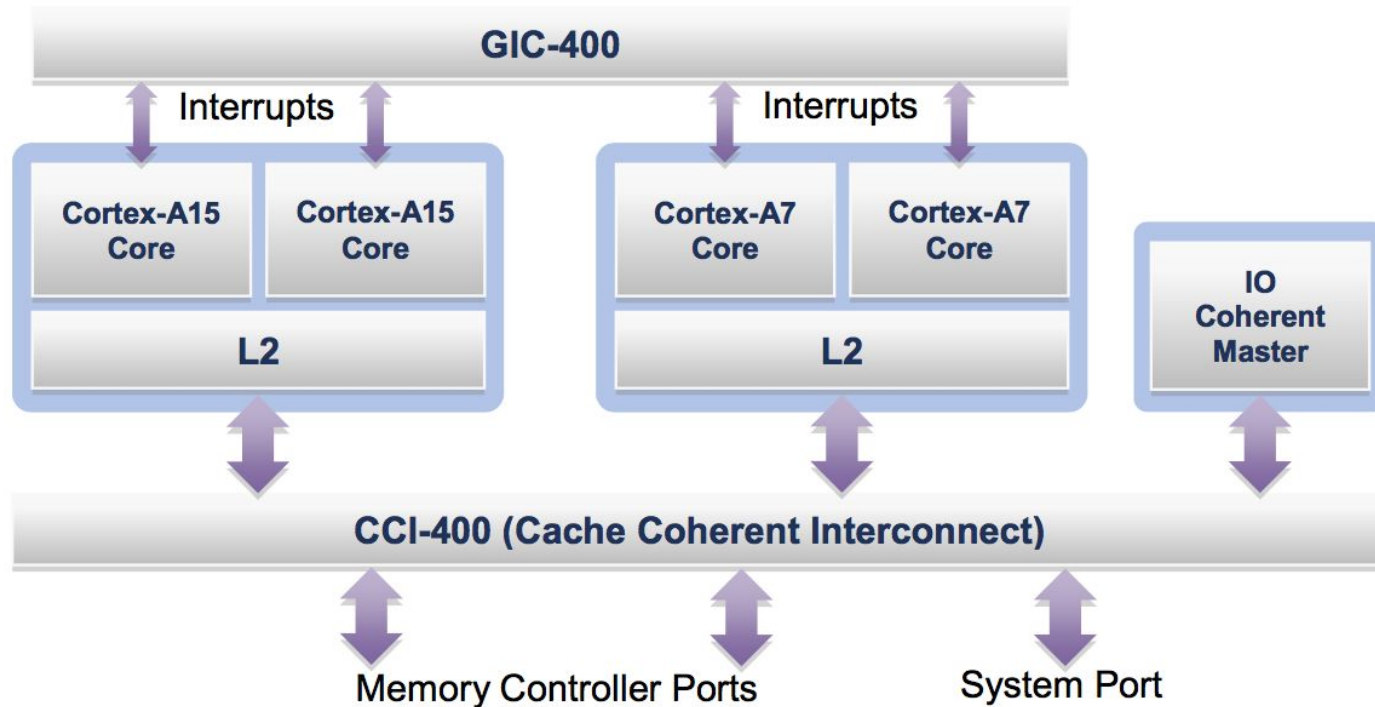
ARM Cortex-A7 (LITTLE) vs Cortex-A15 (Big)

- In-order
- Non-symmetric Dual Issue
- Processor pipeline between 8-10 stages

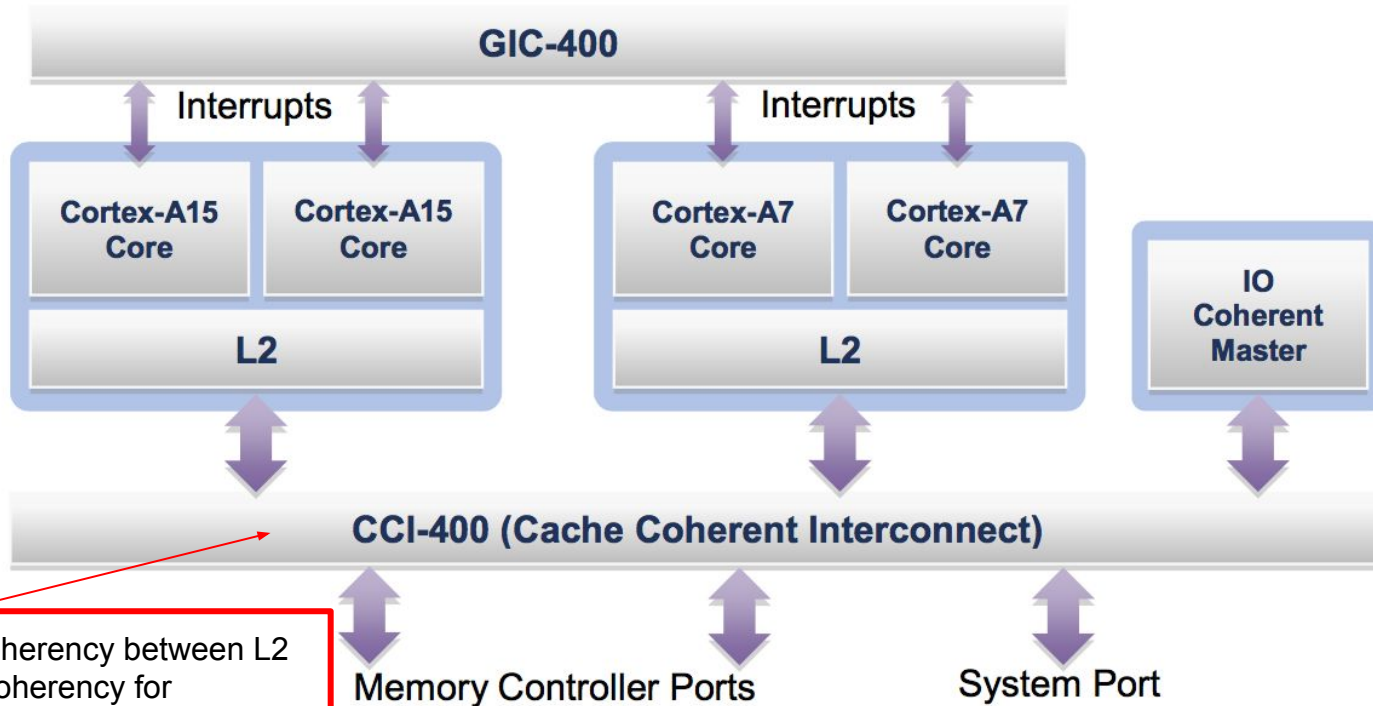


- Out-of-order
- Sustained Triple-issue
- Pipeline between 15-24 stages

Interconnection



Interconnection

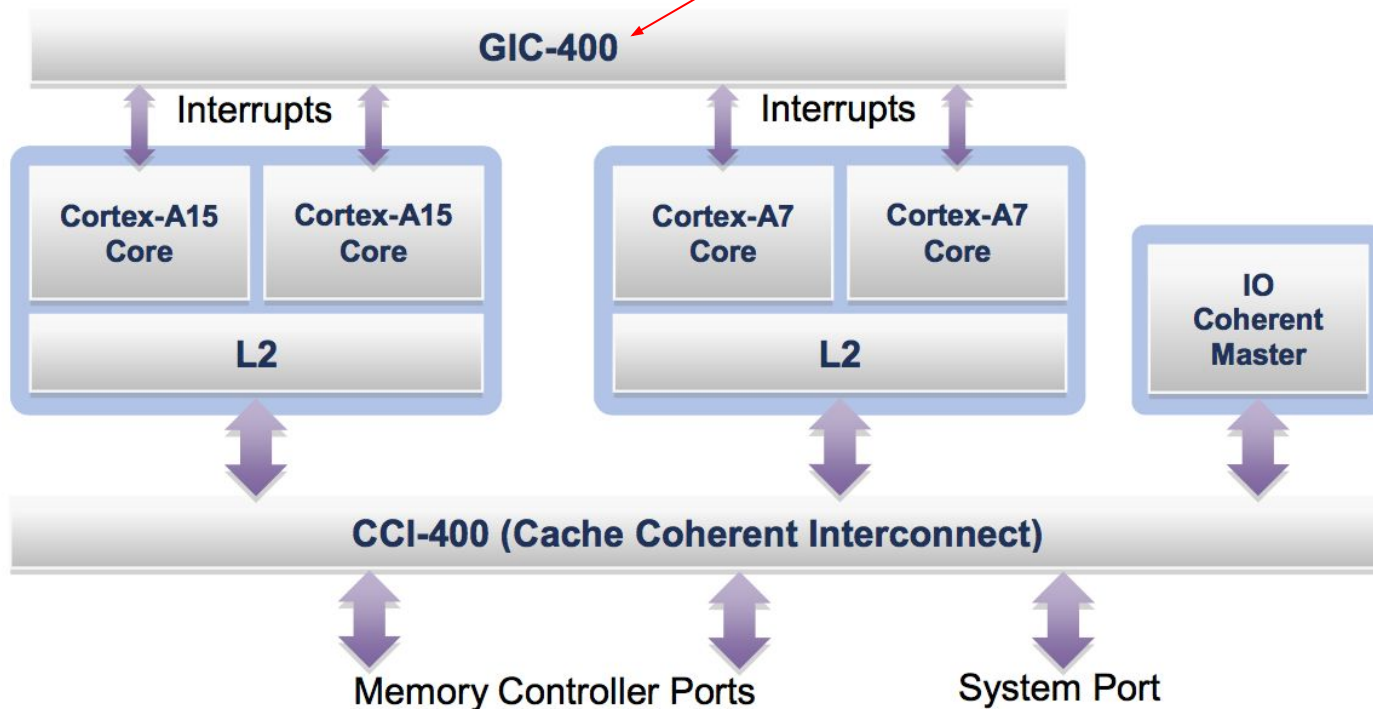


Facilitates full-coherency between L2 caches and IO coherency for components such as GPU.

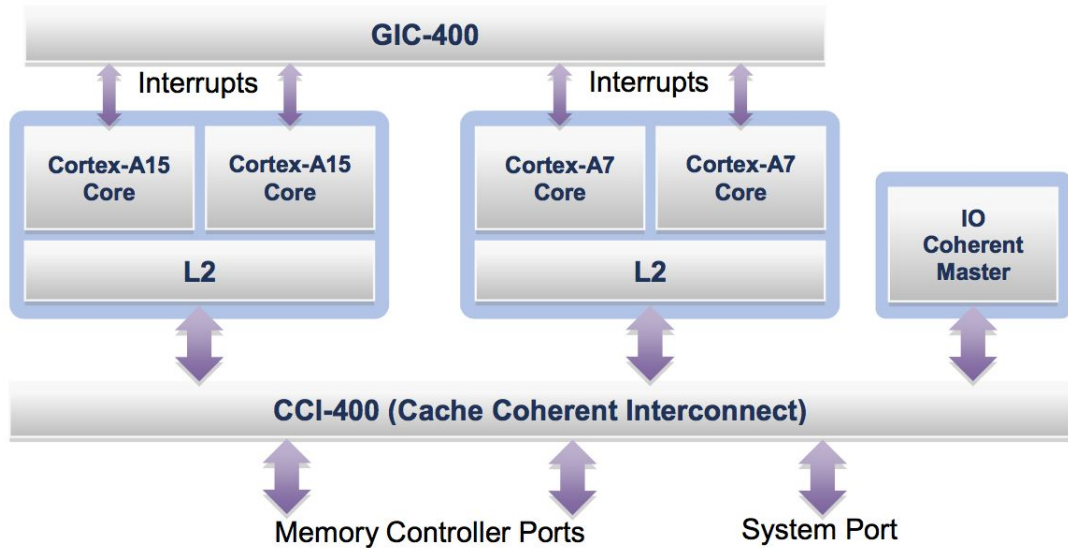
Interconnection

Distribute up to 480-interrupts.

GIC-400 allows interrupts to be migrated between any cores in the Cortex-A15 cluster or Cortex-A7 cluster.



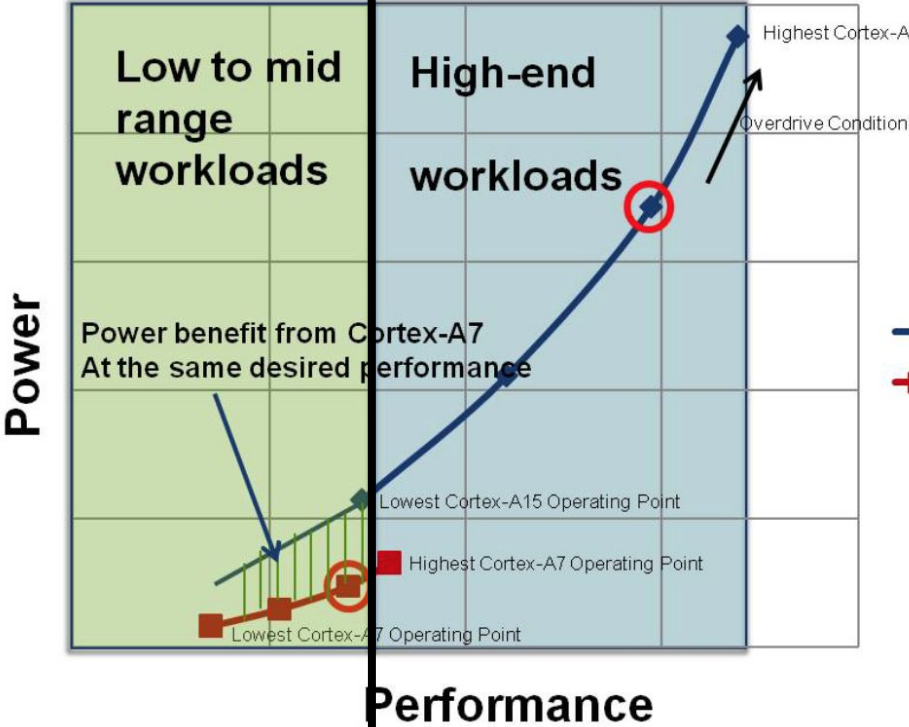
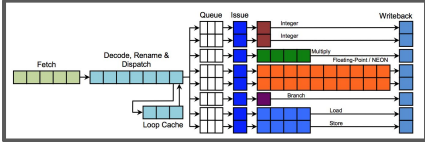
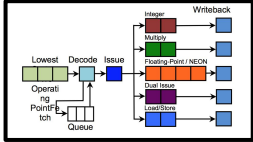
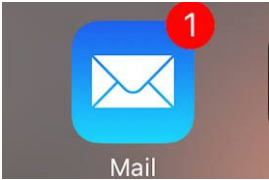
Interconnection



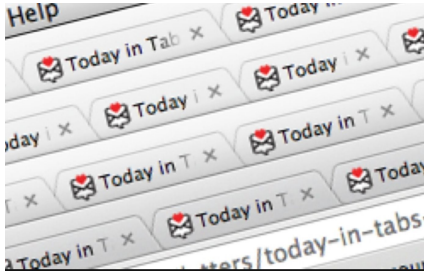
From ARM White Paper:

- Any arrangements of Big.LITTLE is possible.
- There is no feature or optimization for a specific configuration of Big.LITTLE prior from ARM.
- The recommendation is # bigs = # littles (in order to reduce **software complexity** of the task scheduler).

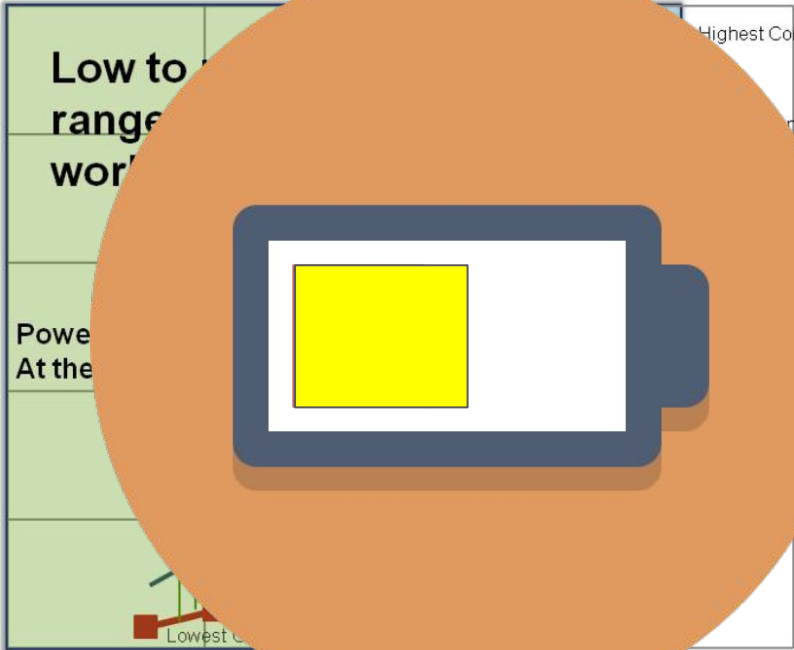
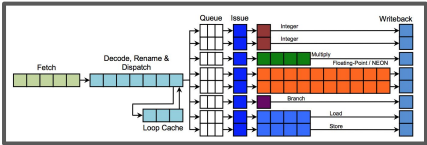
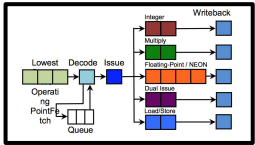
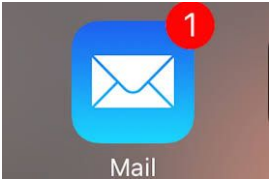
Usage



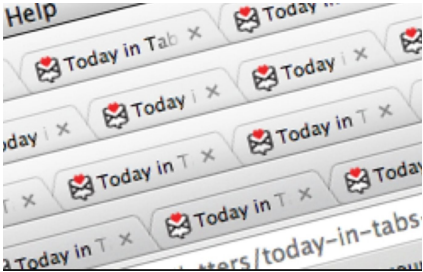
- Cortex-A15
- Cortex-A7



Usage



Cortex-A15
Cortex-A7



Task Migration

Migration between a little to a big core or vice-versa takes less than **20.000 cycles**.



In a 1Ghz processors, it would take **20 us (micro-seconds)**.

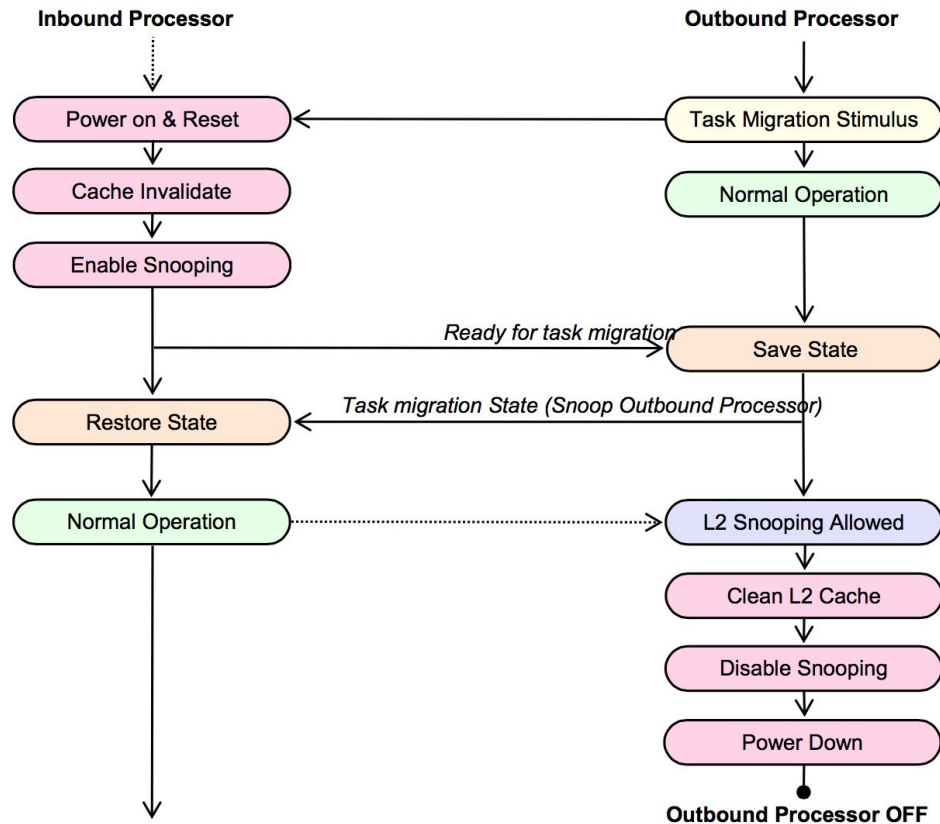


In practice, it goes around 30us

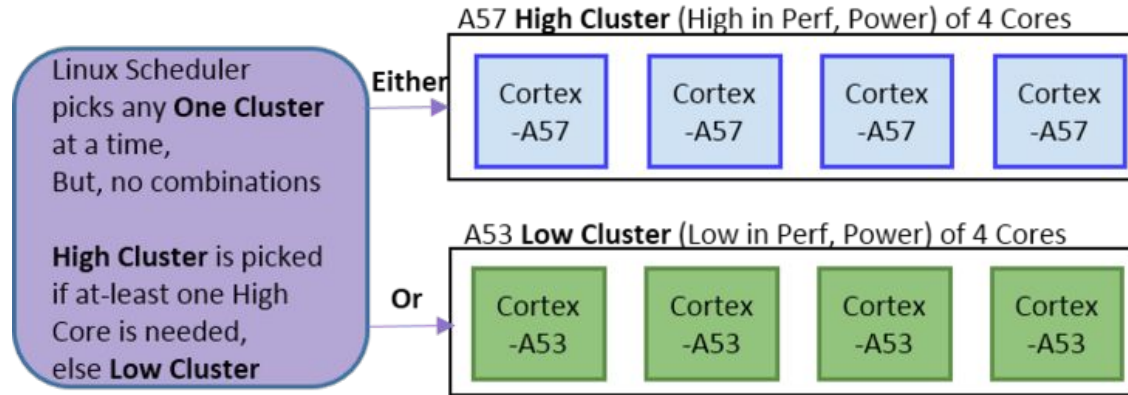


This is low:

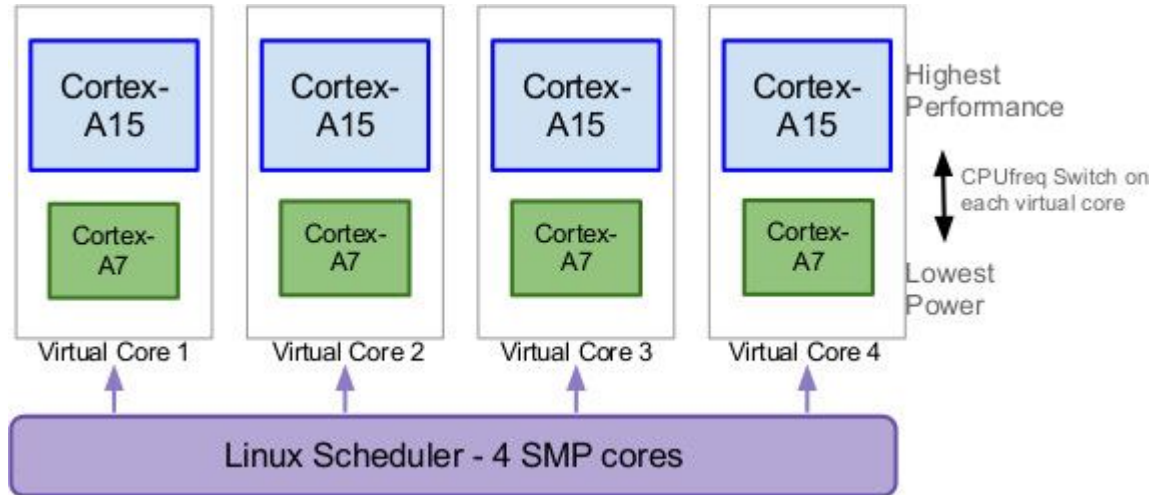
- DVFS I would wait around 100/200 microseconds.
- Any consistent load monitoring takes at least 1 millisecond.



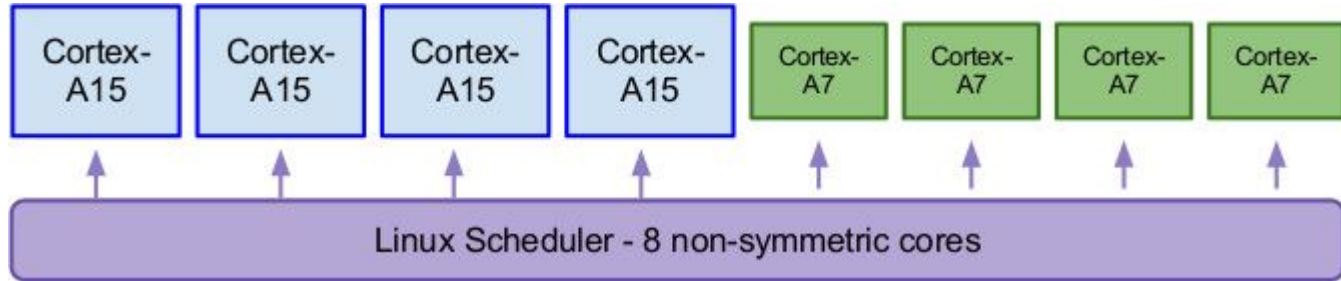
Big.LITTLE arrangement: clustered switching



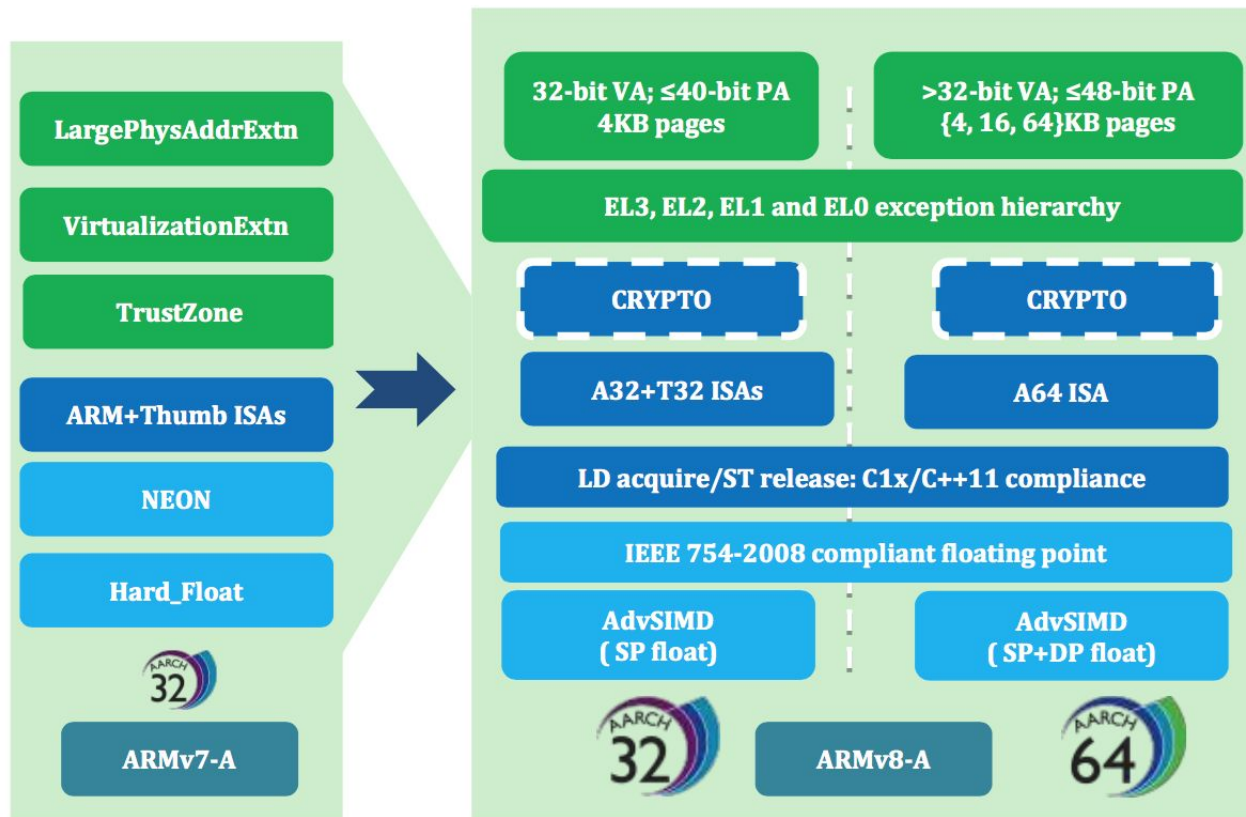
Big.LITTLE arrangement: in-kernel switcher



Big.LITTLE arrangement: heterogeneous multi-processing

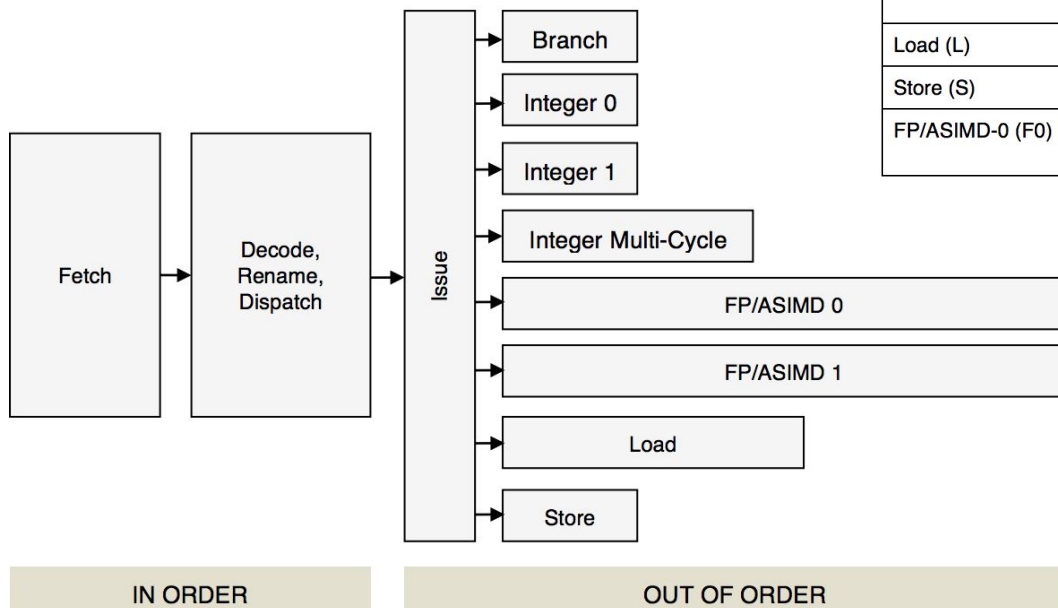


Big.LITTLE on armv8



Big.LITTLE on armv8

A-57



Pipeline (mnemonic)	Supported functionality
Branch (B)	Branch μ ops
Integer 0/1 (I)	Integer ALU μ ops
Multi-cycle (M)	Integer shift-ALU, multiply, divide, CRC and sum-of-absolute-differences μ op μ ops
Load (L)	Load and register transfer μ ops
Store (S)	Store and special memory μ ops
FP/ASIMD-0 (F0)	ASIMD ALU, ASIMD misc, FP misc, FP add, FP multiply, ASIMD integer multiply, FP divide μ ops, crypto μ ops

Big.LITTLE on armv8

A-53

- Versatile, can be paired with any Armv8.0 core in a big.LITTLE configuration, including Cortex-A57, Cortex-A72, other Cortex-A53, and Cortex-A35 processors.

This processor can also be implemented in an Arm big.LITTLE configuration.

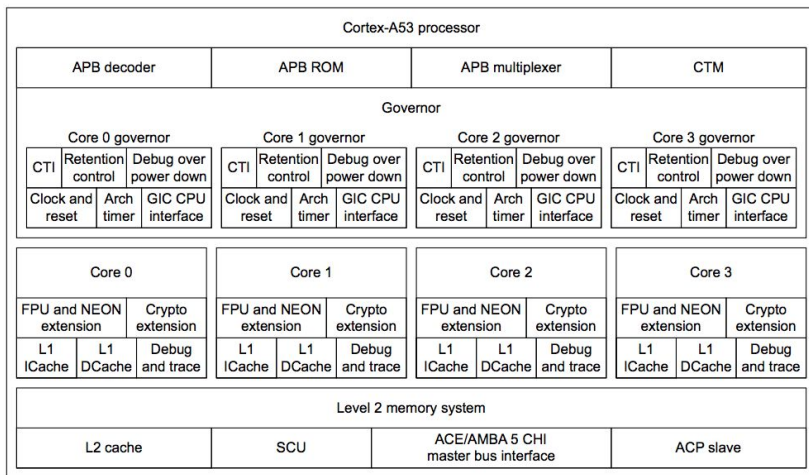
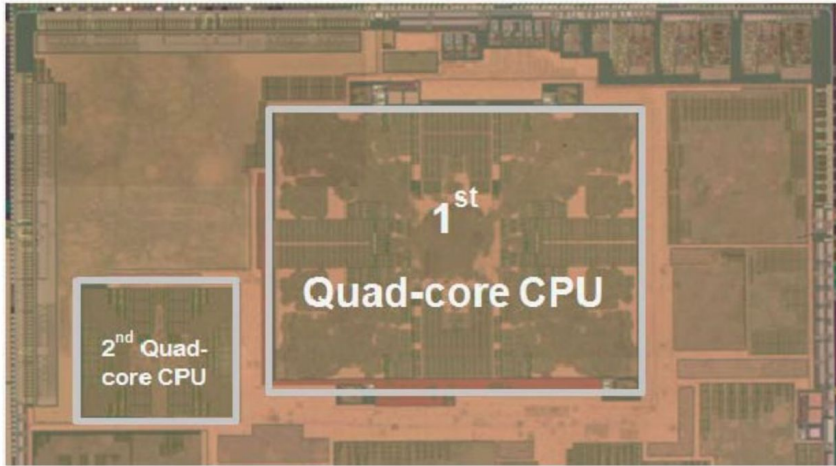


Figure 2-1 Cortex-A53 processor block diagram














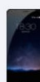











Applications



Inside the first Samsung Exynos Octacore chip

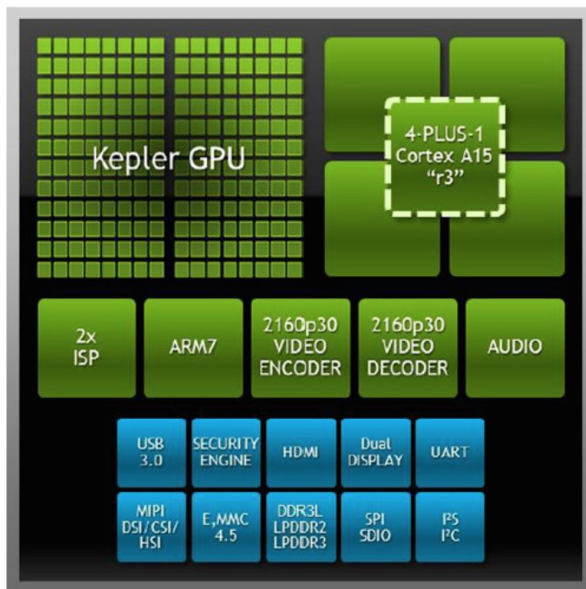
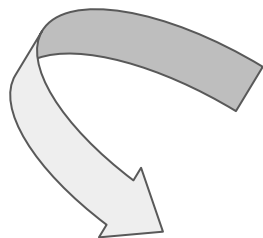
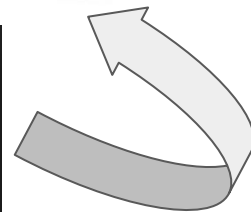
Smartphones with Samsung Exynos 7 Octa 7420 processor

List of smartphones that works with Samsung Exynos 7 Octa 7420 processor inside.

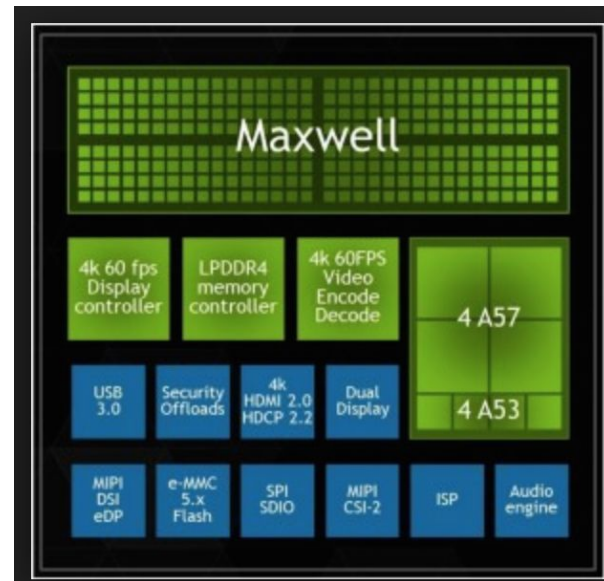
 Samsung Galaxy S6 Edge 32GB €388	 Samsung Galaxy S6 Edge 64GB €428	 Samsung Galaxy S6 Edge 128GB	 Samsung Galaxy S6 128GB
 Samsung Galaxy S6 64GB	 Samsung Galaxy S6 32GB €262	 Samsung Galaxy S6 32GB Dual	 Samsung Galaxy S6 edge+ 64GB G928F EU
 Samsung Galaxy S6 edge+ 32GB G928F EU	 Samsung Galaxy S6 edge+ 32GB Dual SIM G9287	 Samsung Galaxy Note 5 N920C 32GB €410	 Samsung Galaxy Note 5 N920I 32GB
 Samsung Galaxy Note 5 N9200 Dual SIM €484	 Meizu Pro 5 3GB 32GB €375	 Meizu Pro 5 4GB 64GB €431	 Samsung Galaxy Note 5 N9208 Dual SIM
 Meizu Pro 5 3GB 64GB	 Meizu Pro 5 4GB 32GB	 Samsung Galaxy Note 5 N920I 64GB	 Meizu Pro 5 4GB 64GB Internati...
 Samsung Galaxy S6 edge+ 32GB G928C €464	 Samsung Galaxy Note 5 N920I 128GB	 Samsung Galaxy S6 Active	 Samsung Galaxy Note 5 N920C 64GB €410
 Samsung Galaxy A8 (2016) A810F			

<https://www.kimovil.com/en/list-smartphones-by-processor/samsung-exynos-7-octa-7420>

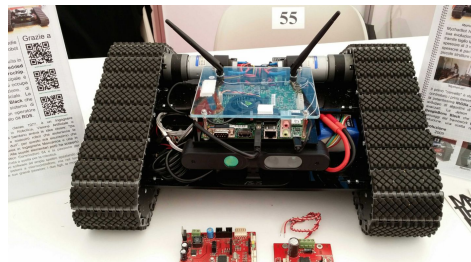
Applications



NVIDIA Tegra TK1 - 4 Plus 1



NVIDIA Tegra X1 - 4-4



A bit of History

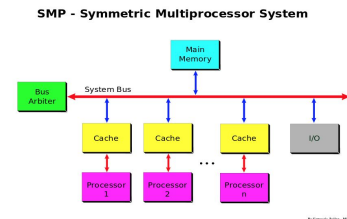
Homogeneous computing: adequate performance for many applications in the past.

- More parallelization
- Heterogeneous consumption of resources from applications
- Opportunity: initially for thinking on performance

Heterogeneous Computing

“is the well-orchestrated and coordinated effective use of a **suite of diverse** high-performance **machines** (including parallel machines) to provide superspeed processing for **computationally demanding tasks with diverse computing needs**. An HC system includes heterogeneous machines, high-speed networks, interface.” (1993)

“refers to the use of **different processing cores** to maximize performance” (2010)

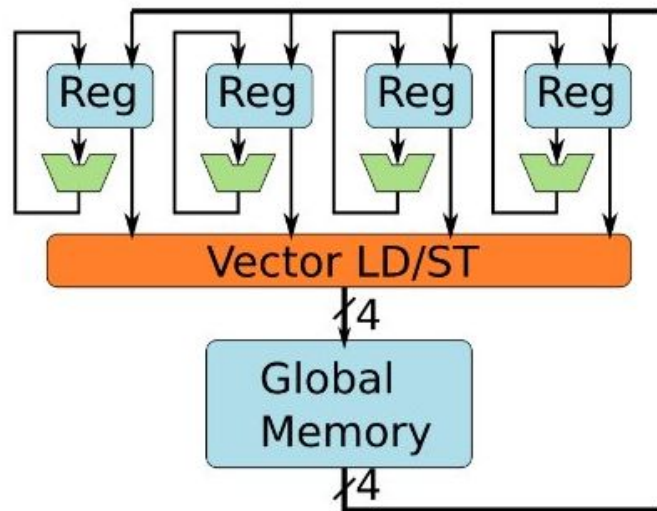


Accelerators

- CPU sometimes is not good enough!
 - SSEs
 - GPUs
 - FPGAs

Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs
 - FPGAs



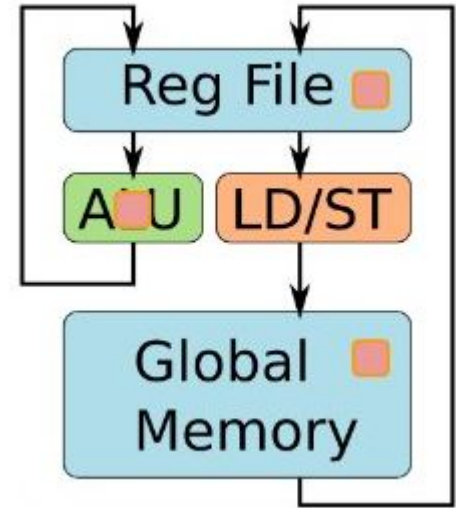
Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs
 - FPGAs
- CPU example

```
int A[2][4];  
for(i=0;i<2;i++){  
    for(j=0;j<4;j++){  
        A[i][j]++;  
    }  
}
```

Assembly
code of
inner-loop

```
lw r0, 4(r1)  
addi r0, r0, 1  
sw r0, 4(r1)
```

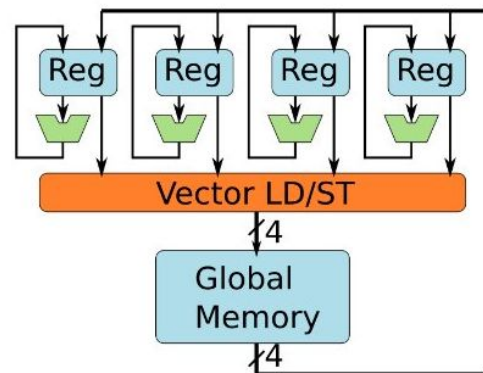


Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs
 - FPGAs
- SSE example - 4 ALUs

```
int A[2][4];  
for(i=0;i<2;i++){  
    for(j=0;j<4;j++){  
        A[i][j]++;  
    }  
}
```

```
int A[2][4];  
for(i=0;i<2;i++){  
    movups xmm0, [ &A[i][0] ] // Load  
    addps xmm0, xmm1 // add 1  
    movups [ &A[i][0] ], xmm0 // store  
}
```



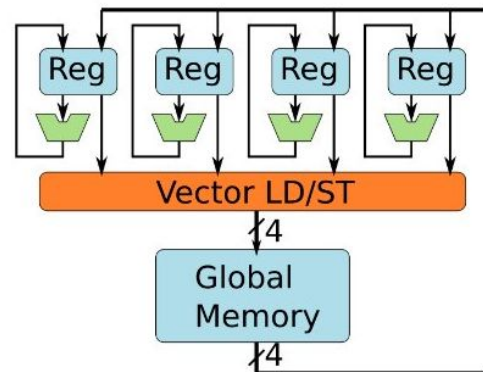
Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs
 - FPGAs
- SSE
 - X86 Pentiums - AMD Athlon XP

```
int A[2][4];  
for(i=0;i<2;i++){  
    for(j=0;j<4;j++){  
        A[i][j]++;  
    }  
}
```

→

```
int A[2][4];  
for(i=0;i<2;i++){  
    movups xmm0, [ &A[i][0] ] // Load  
    addps xmm0, xmm1 // add 1  
    movups [ &A[i][0] ], xmm0 // store  
}
```

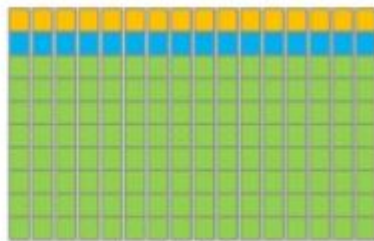


Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs (Graphics Processing Unit)
 - FPGAs



CPU

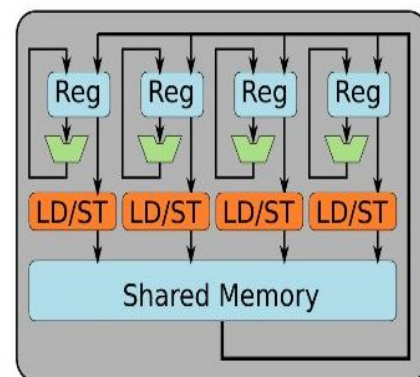
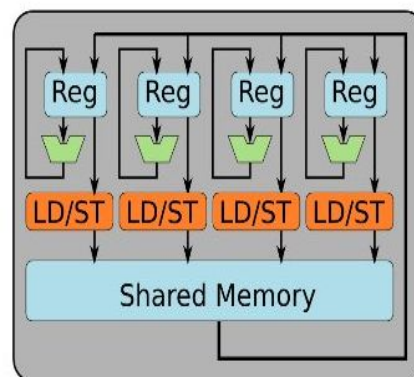
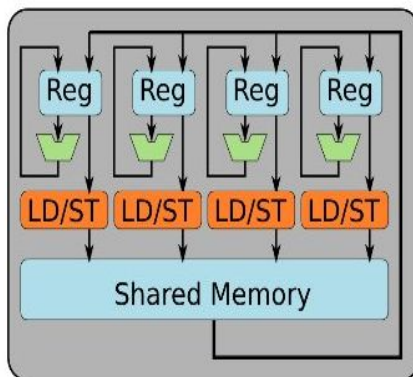
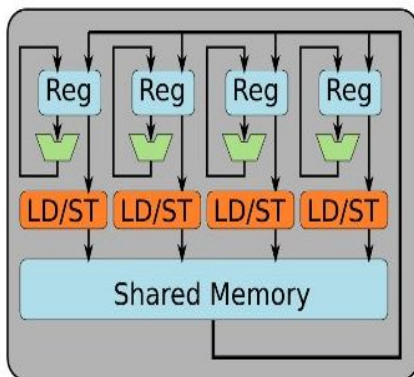


GPU



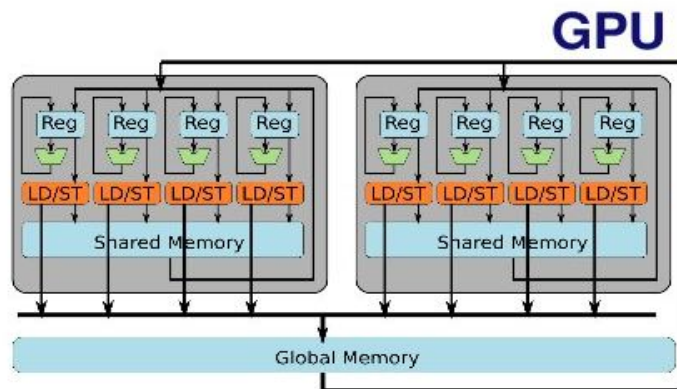
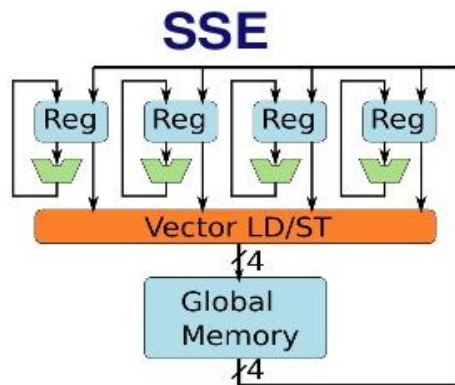
Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - **GPUs (Graphics Processing Unit)**
 - FPGAs
- GPU contains multiple SIMD cores.



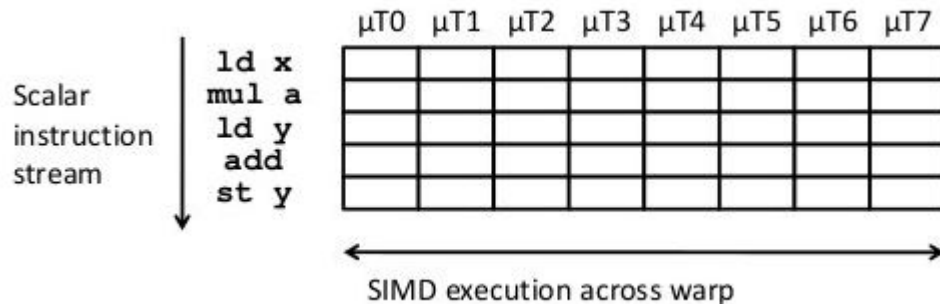
Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs (Graphics Processing Unit)
 - FPGAs
- Shared memory



Accelerators

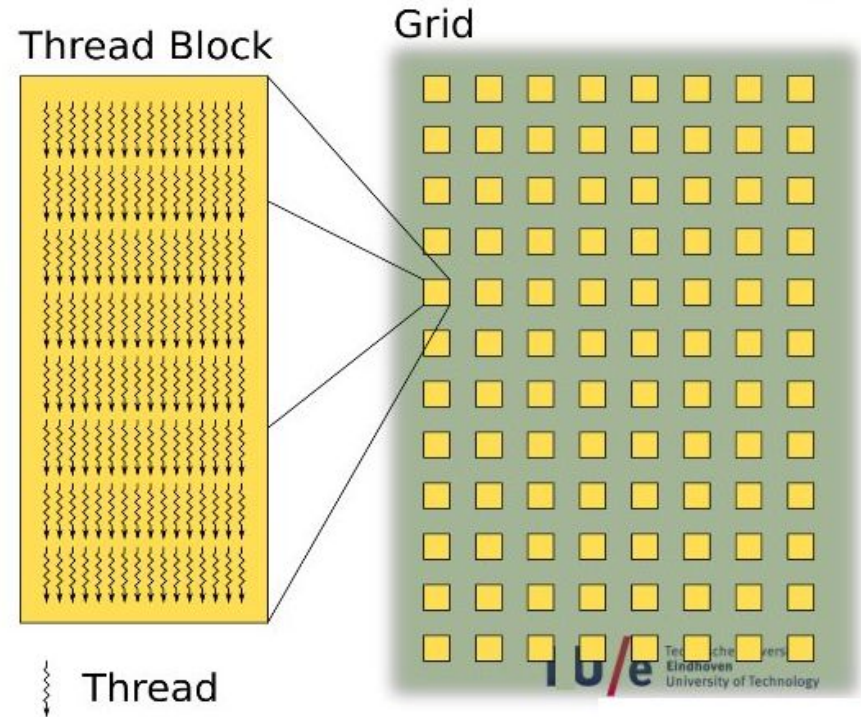
- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs (Graphics Processing Unit)
 - FPGAs
- SIMT (Single Instruction Multiple Threads)



SIMD execution across warp

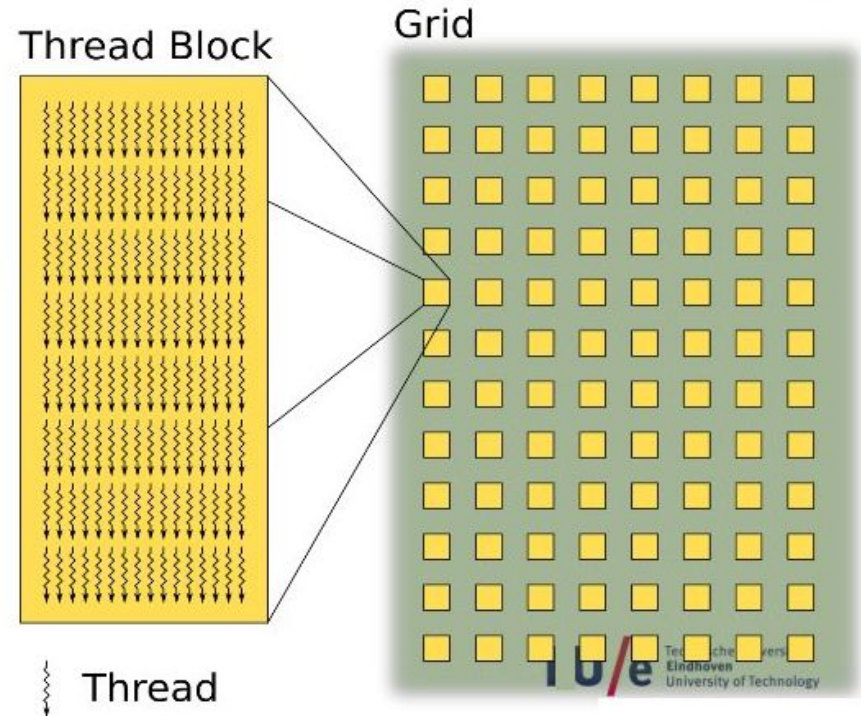
Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - **GPUs (Graphics Processing Unit)**
 - FPGAs
- Thread Hierarchy
 - Grid contains Blocks
 - Blocks contains Threads



Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - **GPUs (Graphics Processing Unit)**
 - FPGAs
- Thread Hierarchy
 - Grid contains Blocks
 - Blocks contains Threads
 - NVIDIA
 - 32 threads = *warp*

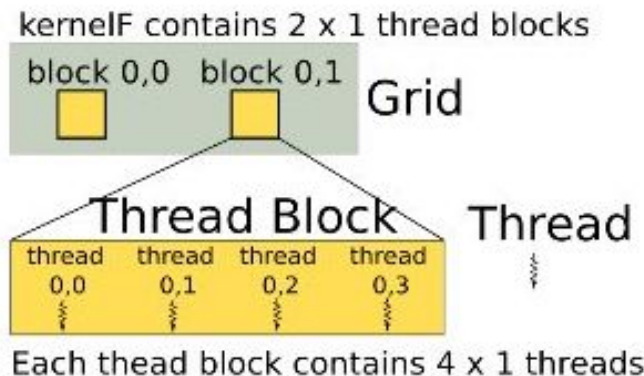


Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs (Graphics Processing Unit)
 - FPGAs
- Thread Hierarchy
 - Grid contains Blocks
 - Blocks contains Threads
 - NVIDIA
 - 32 threads = *warp*

```
int A[2][4];  
kernelF<<<(2,1),(4,1)>>>(A); // define threads
```

```
__device__ kernelF(A){ // all threads run same kernel  
    i = blockIdx.x;      // each thread block has its id  
    j = threadIdx.x;      // each thread has its id  
    A[i][j]++;           // each thread has a different i and j  
}
```



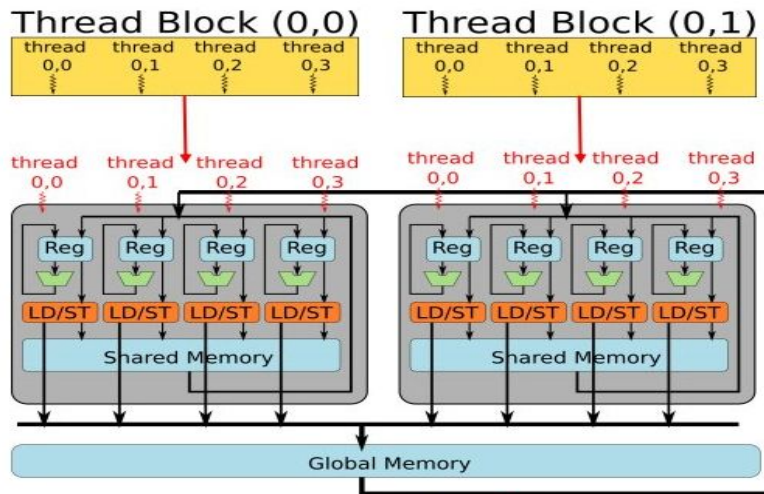
Accelerators

- CPU sometimes is not good enough!
 - SSEs (Stream SIMD Extensions)
 - GPUs (Graphics Processing Unit)
 - FPGAs
- Thread Hierarchy
 - Grid contains Blocks
 - Blocks contains Threads
 - NVIDIA
 - 32 threads = *warp*

```
int A[2][4];  
kernelF<<<(2,1),(4,1)>>>(A); // define threads
```

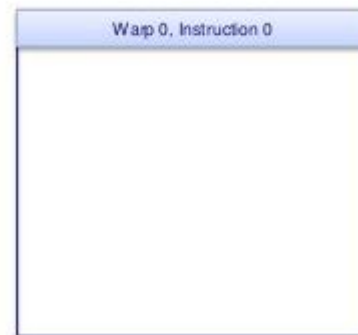
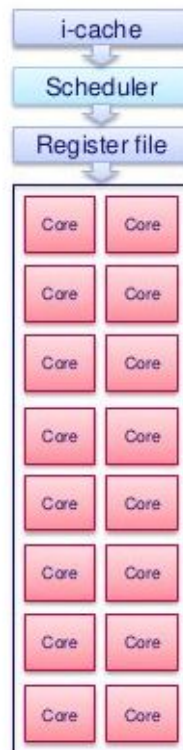
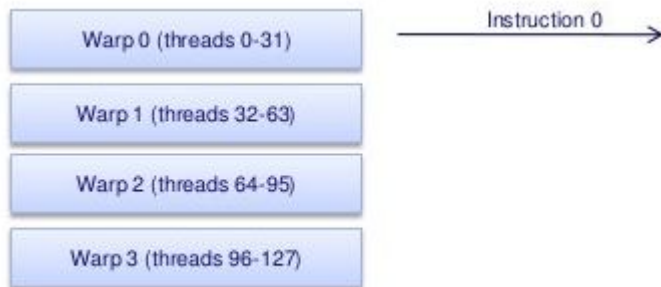
```
__device__ kernelF(A){ // all threads run same kernel  
    i = blockIdx.x; // each thread block has its id  
    j = threadIdx.x; // each thread has its id  
    A[i][j]++; // each thread has a different i and j  
}
```

TU/e



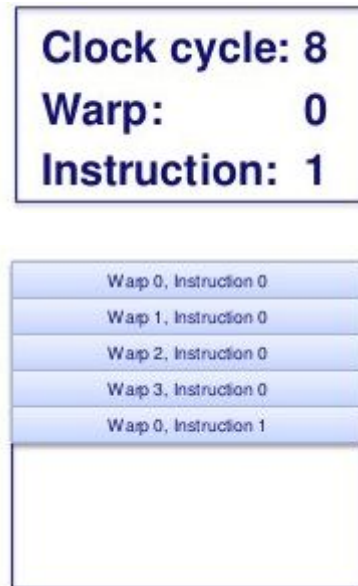
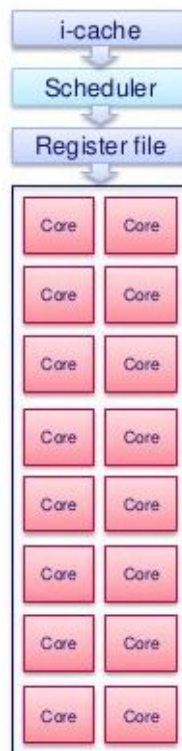
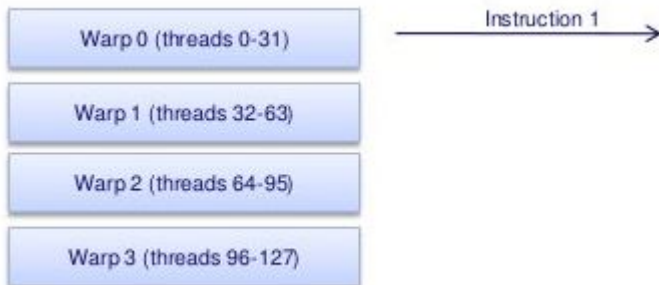
Accelerators

- Warp Scheduler
 - Example
 - 128 threads
 - 4 warps
 - 16 SIMD cores (Stream Multiprocessor)



Accelerators

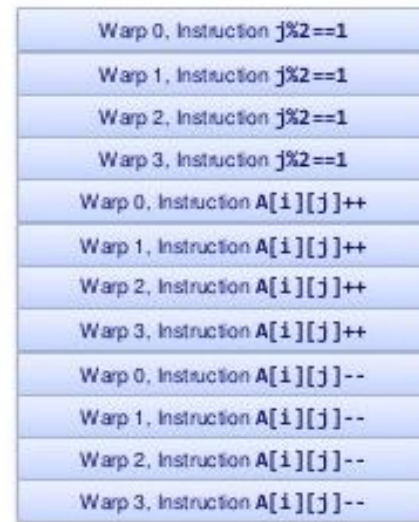
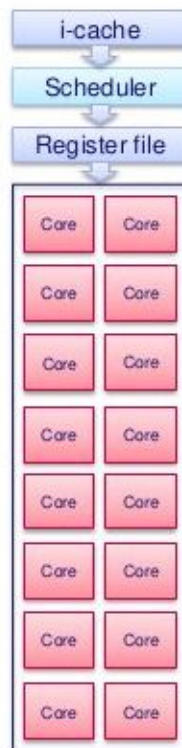
- Warp Scheduler
 - Example
 - 128 threads
 - 4 warps
 - 16 SIMD cores (Stream Multiprocessor)



Accelerators

- Warp Scheduler
 - Example
 - 128 threads
 - 4 warps
 - 16 SIMD cores (Stream Multiprocessor)

```
__device__ kernelF(A){  
    i = blockIdx.x;  
    j = threadIdx.x;  
    if(j%2 == 1)  
        A[i][j]++;  
    else  
        A[i][j]--;  
}
```



Accelerators

- Warp Scheduler

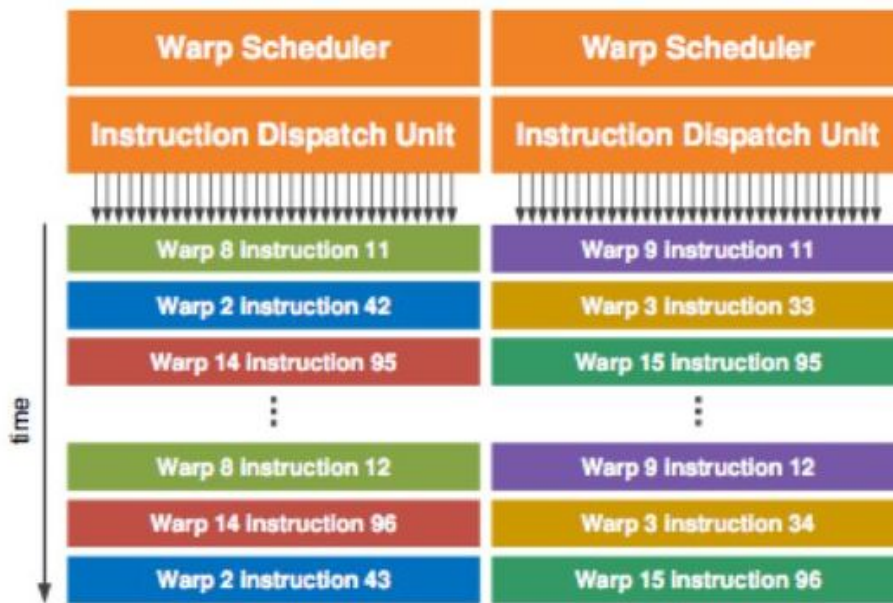


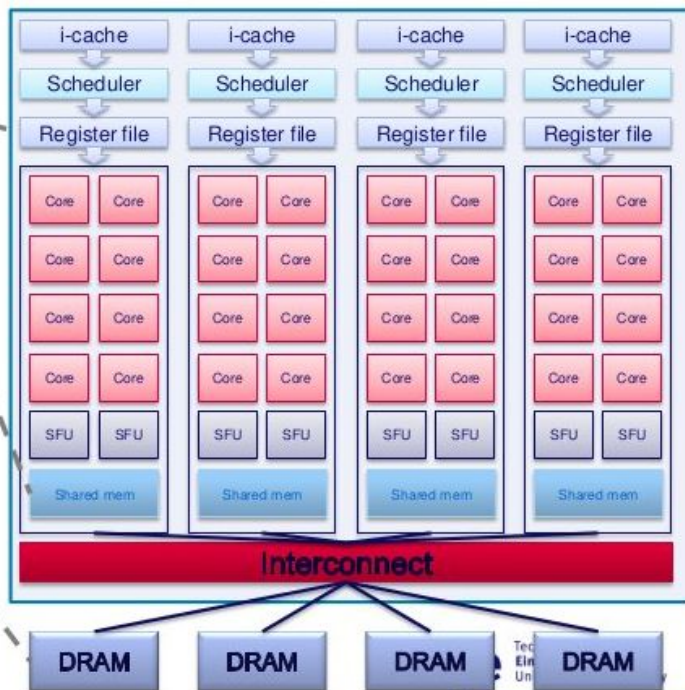
Image Credit: NVIDIA

http://www.nvidia.com/content/PDF/fermi_white_papers/NVIDIA_Fermi_Compute_Architecture_Whitepaper.pdf

Accelerators

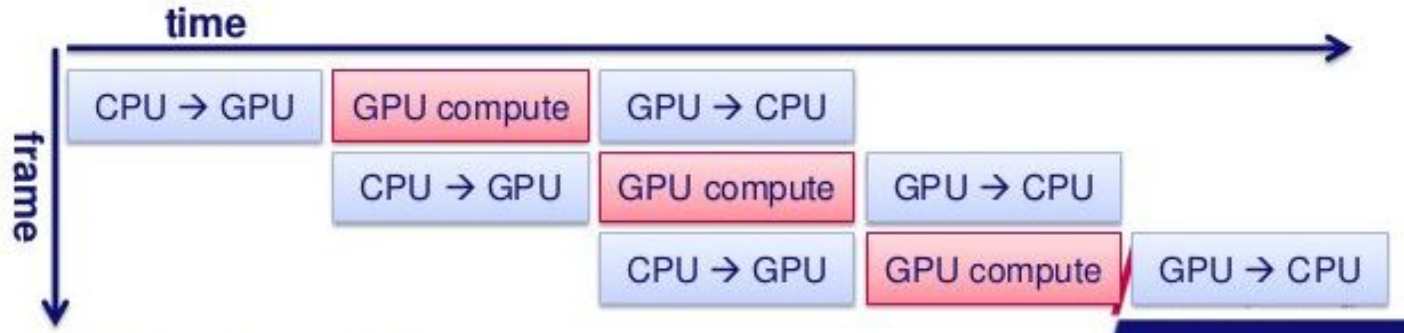
- Memory Hierarchy

- **Register**
 - Per thread
- **Shared memory**
 - Per thread block
- **Global memory**
 - Per kernel



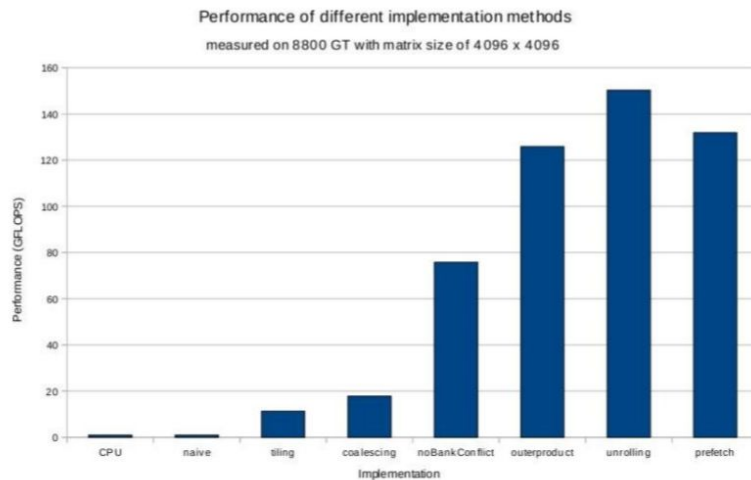
Accelerators

- Memory Hierarchy
- Takes long time to transmit data from CPUs to GPGPUs



Accelerators

- Memory Hierarchy
- Takes long time to transmit data from CPUs to GPGPUs
- Matrix Multiplication



Accelerators

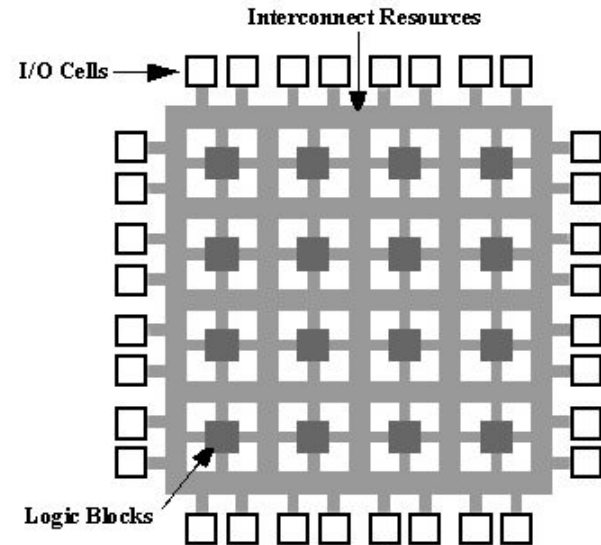
Nvidia Fermi GF100 GPU



Image Credit: NVIDIA [Wittenbrink, Kilgariff, and Prabhu, Hot Chips 2010]

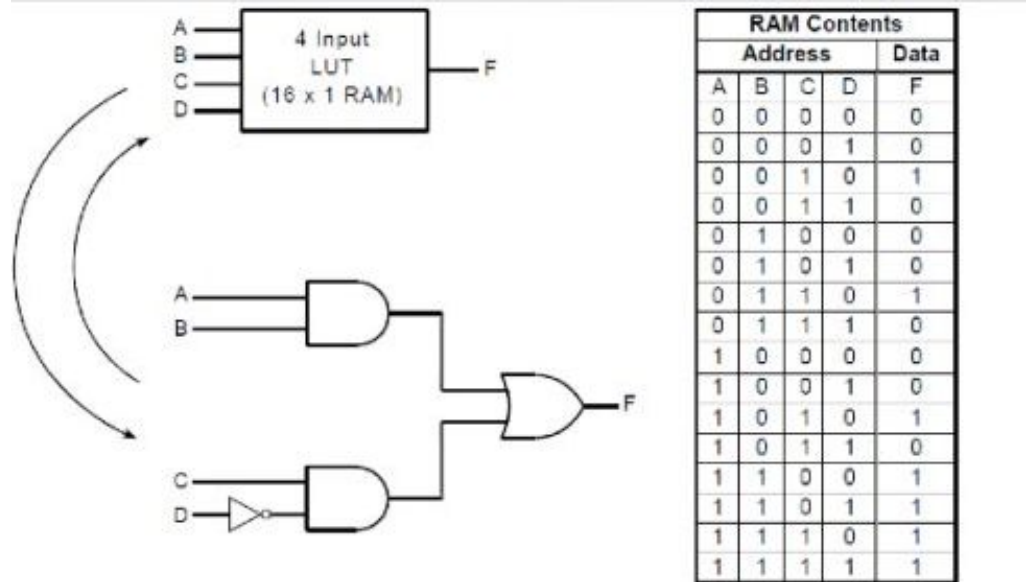
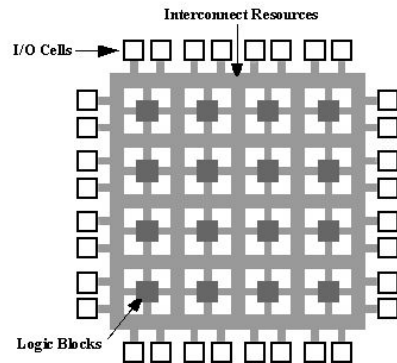
Accelerators

- CPU sometimes is not good enough!
 - SSEs
 - GPUs
 - **FPGAs (Field-Programmable Gate Array)**



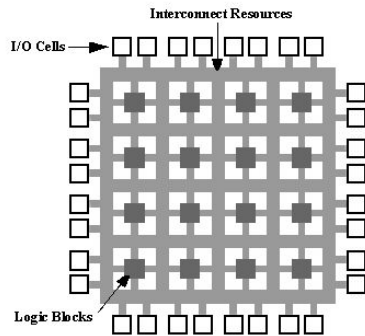
Accelerators

- CPU sometimes is not good enough!
 - SSEs
 - GPUs
 - **FPGAs**
- Logic as Memories
 - LUT (LookUp Table)

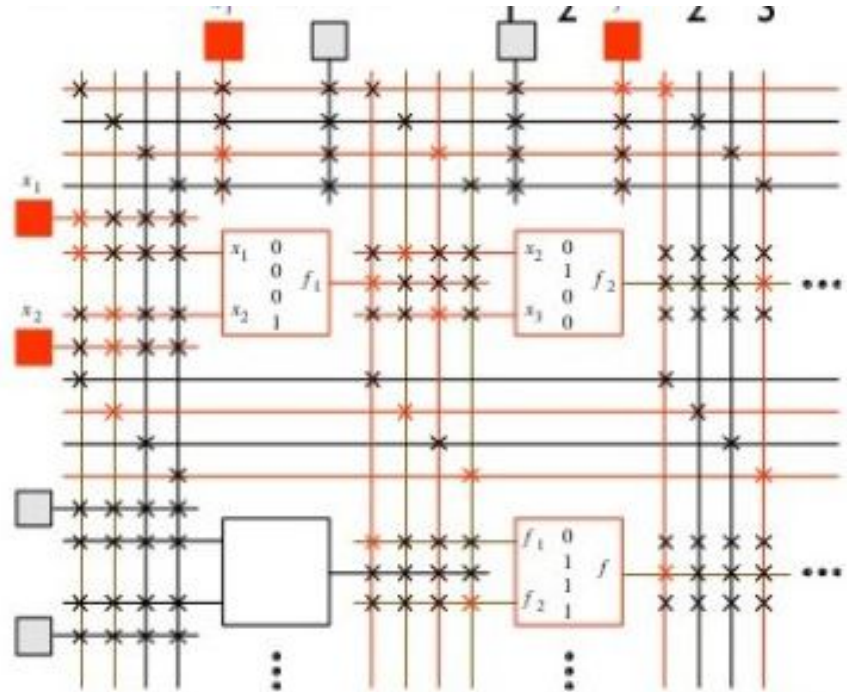


Accelerators

- CPU sometimes is not good enough!
 - SSEs
 - GPUs
 - **FPGAs**
- Interconnection

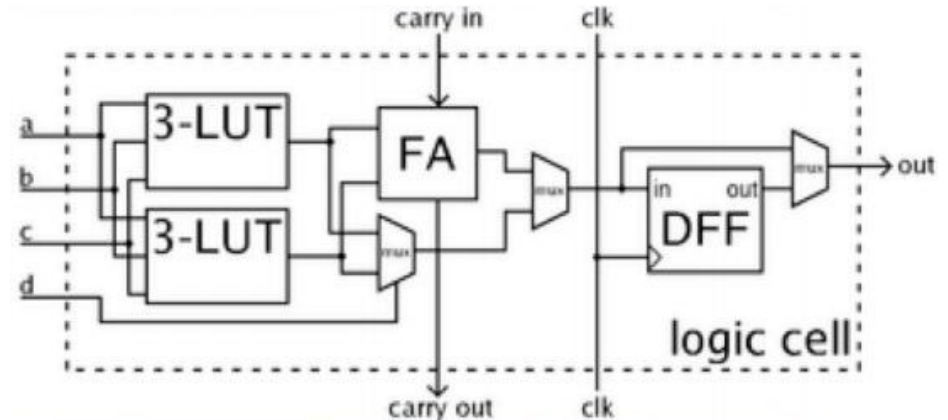
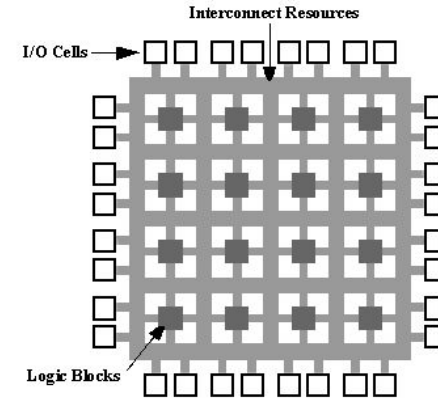
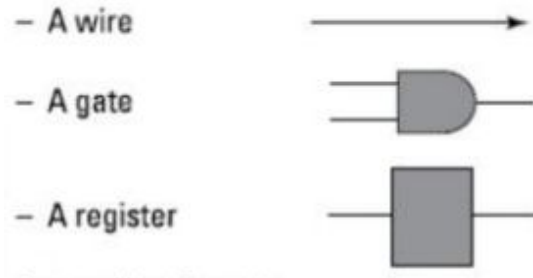


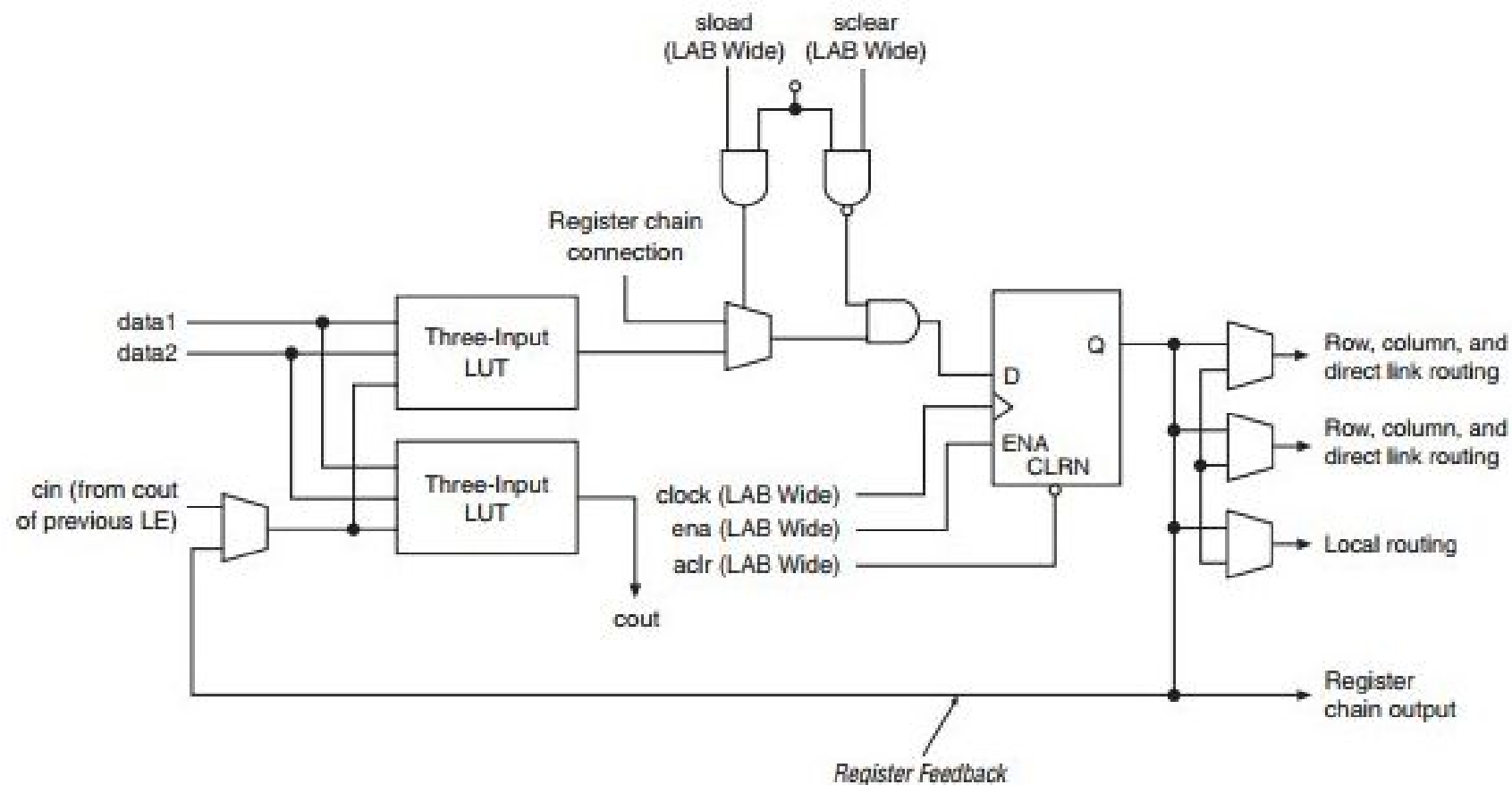
- $f_1 = x_1 x_2$
- $f_2 = x_2' x_3$
- $f = f_1 + f_2$



Accelerators

- CPU sometimes is not good enough!
 - SSEs
 - GPUs
 - **FPGAs**
- Logical Blocks





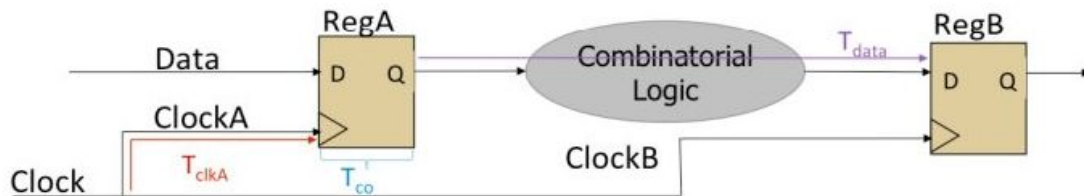
Accelerators

- CPU sometimes is not good enough!
 - SSEs
 - GPUs
 - **FPGAs**
 - Programming
 - Analysis & Synthesis (From HDL)
 - Fitter
 - Assembler
 - Timing Analysis



Accelerators

- CPU sometimes is not good enough!
 - SSEs
 - GPUs
 - **FPGAs**
 - Programming
 - Analysis & Synthesis (From HDL)
 - Fitter
 - Assembler
 - **Timing Analysis**



Accelerators

- CPU sometimes is not good enough!
 - SSEs
 - GPUs
 - **FPGAs**
 - Custom Blocks
 - DSPs
 - Multipliers
 - Embedded Memories
 - Applications
 - Image Processing
 - DNN (!!)
 - Prototyping
 - SoC FPGAs (CPU + FPGA)



References

Khokhar, Ashfaq A., et al. "Heterogeneous computing: Challenges and opportunities." Computer 26.6 (1993): 18-27.

Brodtkorb, Andre R., et al. "State-of-the-art in heterogeneous computing." Scientific Programming 18.1 (2010): 1-33.

Coursera: Computer Architecture (Princeton - David Wentzlaff); Introduction to FPGA Design for Embedded Systems (University of Colorado Boulder)

Henk Corporaal, Gert-Jan van den Braak (Introduction to GPGPU and GPU-architectures)