

ACCELERATOR DESIGN WITH OPENCL

(ATHENS WEEK 19-24 MARCH, 2018)



PARALLELISM

- Task Parallelism
- Data Parallelism
- Pipeline



TASK PARALLELISM

- A process is broken into different tasks.
- If the tasks are independent they can be run on parallel.
- e.g Car Assembly



DATA PARALLELISM

- Different parts of data can be processed in parallel.
- e.g vector addition

QUIZ

- Get the average value of a vector
Is it a data parallel application ?



PIPELINE

- Parallelism in time.
- each stage is dependent on the previous one.
- e.g Car Assembly, Video Pipeline

QUIZ

- What could be the disadvantages of pipeline ?



QUIZ

- How to calculate the performance of a pipeline ?

PARALLELISM: AMDAHL'S LAW

$$speedup = \frac{1}{S + \frac{P}{N}}$$

- S: Fraction of the application that is serial.
- P: Fraction of the application that is parallelizable.
- N Processor Speedup.



QUIZ

- Suppose that we want to enhance the processor used for web serving. The new processor is 10 times faster on computation in the web serving application than the old processor. Assuming that the original processor is busy with computation 40% of the time and is waiting for I/O 60% of the time, what's the overall speedup gained by incorporating enhancement?

RECAP: COMPUTER ARCHITECTURE

- Processor
- MMU
- Cache
- Main Memory (DDR)

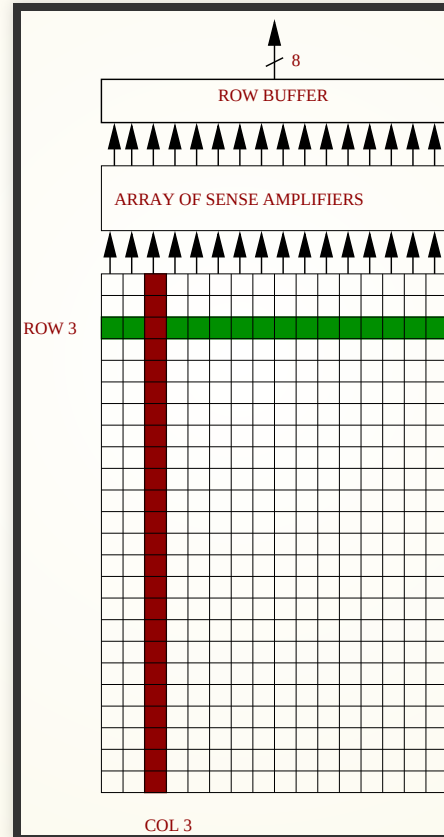


RECAP: COMPUTER ARCHITECTURE

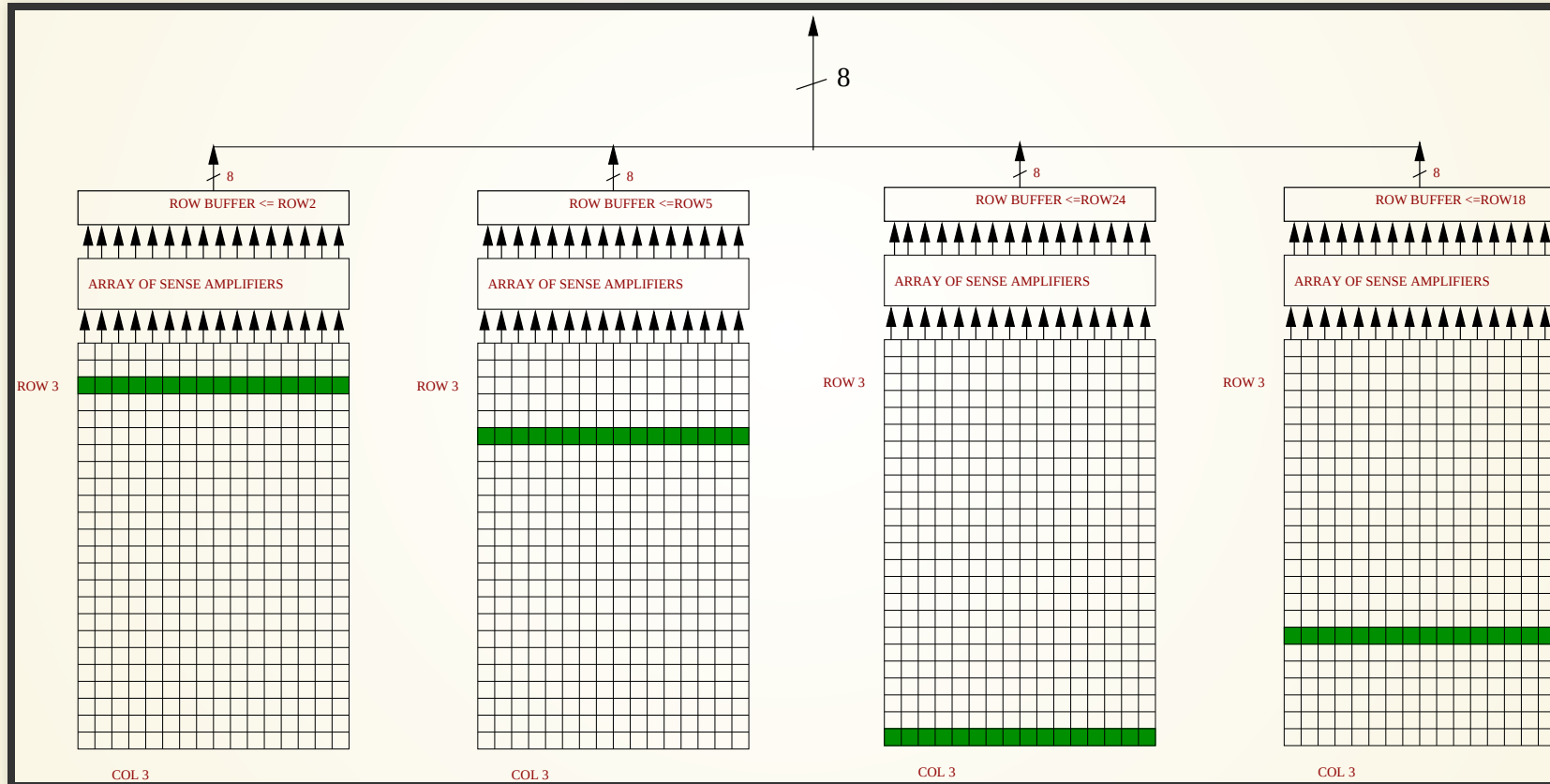
- Main Memory (DDR)
- Cache
- MMU
- Processor



DRAM OPERATION: A SINGLE DRAM BANK



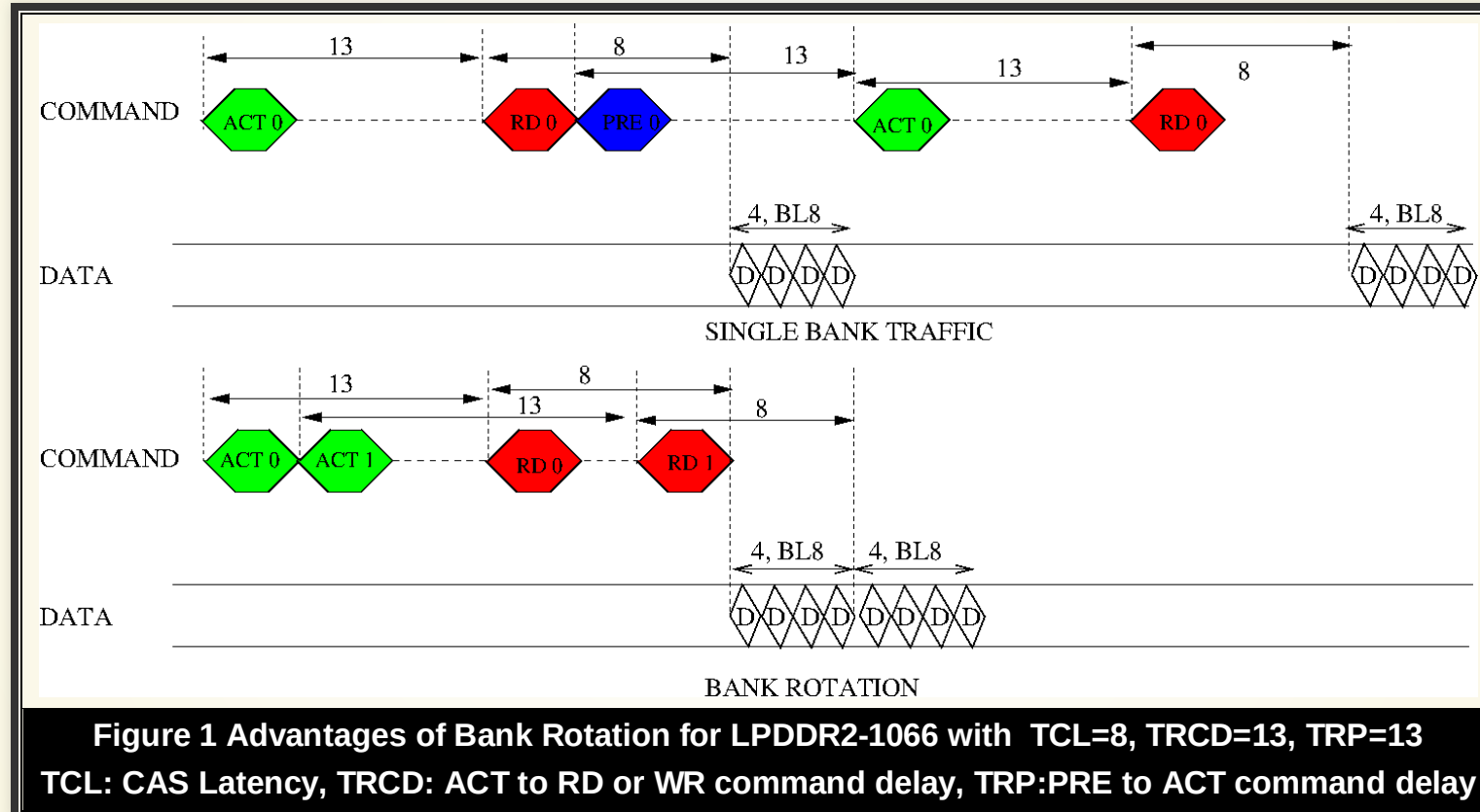
DRAM OPERATION: BANKS



DRAM OPERATION: DRAM OPERATION

- READ: Activate (open the row)-> Read -> Precharge (close).
- WRITE: Activate (open the row)-> Write -> Precharge (close).
- REFRESH: READ-> WRITE back.

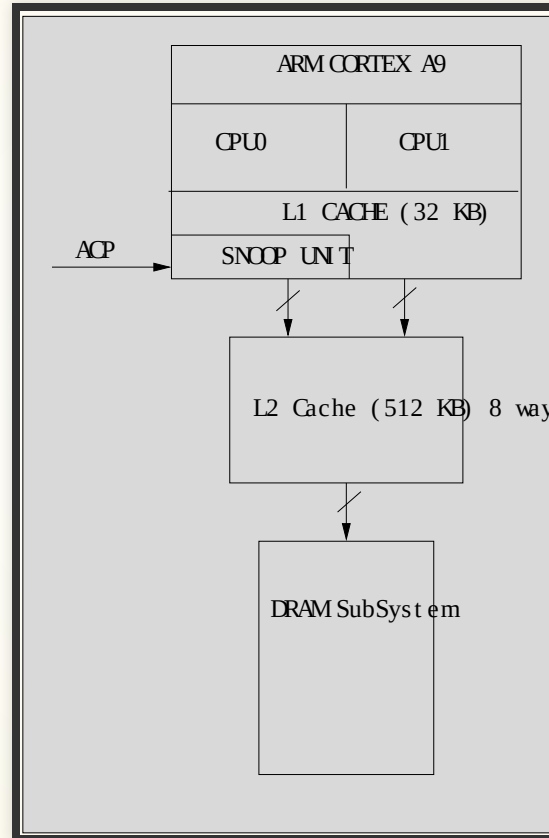
DRAM TRAFFIC TCL,TRCD,TRP



DRAM TRAFFIC

- DRAM is the main performance bottleneck in a SoC.
- DRAM response can come out of order, has high initial latency.

RECAP: CACHE



RECAP: CACHE: CACHE TERMINOLOGY

- Memory contains up-to-date data, and cache has a copy (cache line): CLEAN
- Cache has up-to-date data, and it must be written back to memory: DIRTY
- Memory contains up-to-date data, and cache does not : INVALID
- Memory does not have up-to-date data, cache does not : INVALID

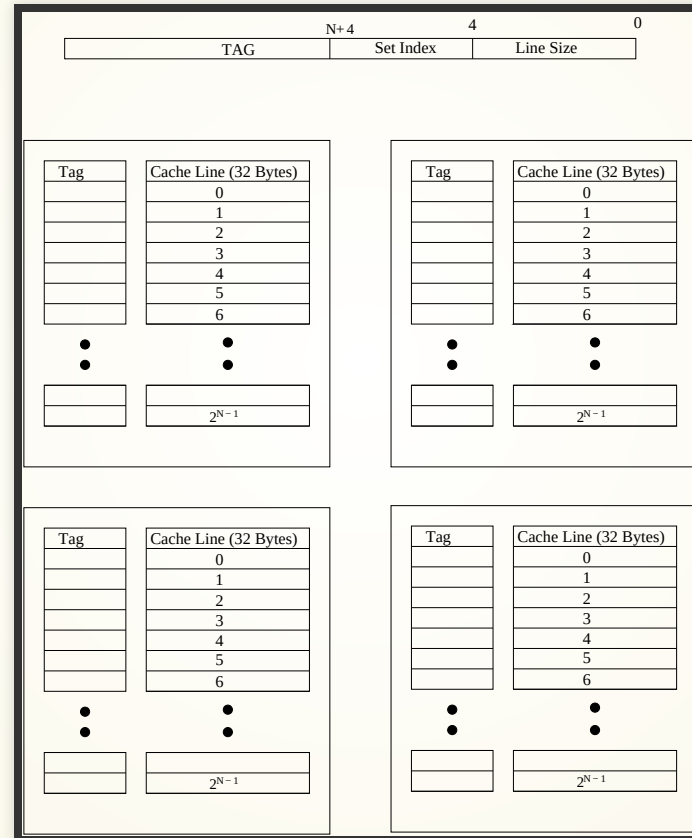


RECAP CACHE TERMINOLOGY

- HIT: Data found in Cache.
- MISS: Data is not in the cache.
- EVICT: A clean cache line is replaced due to a new allocation.



RECAP CACHE ORGANIZATION (4 WAY)



RECAP CACHE POLICIES

- Allocation
 - Write Allocate : On a Write miss replace the cache line.
 - Read Allocate : On a read miss replace the cache line.

RECAP CACHE POLICIES

- Update
 - Write Through : A write updates both the cache and the main memory.
 - Write Back: Write updates the cache only (marked as dirty). Main memory is updated, when the line is evicted, cache is flushed.

RECAP: CACHE COHERENCE

- Case 1. Memory update by another master. Cached copy is out of date.
- Case 2. For write back cache, when master writes to cache, main memory is out of date.



RECAP: CACHE COHERENCE

- Cache Coherency Protocols
 - MEI (Modified, Exclusive, Invalid)
 - MESI (Modified, Exclusive, Shared Invalid)
 - MOESI (Modified, Owned, Exclusive, Shared Invalid)
- Goals
 - Cache to Cache copy of clean data.
 - Cache to Cache move of Dirty data without accessing external memory.

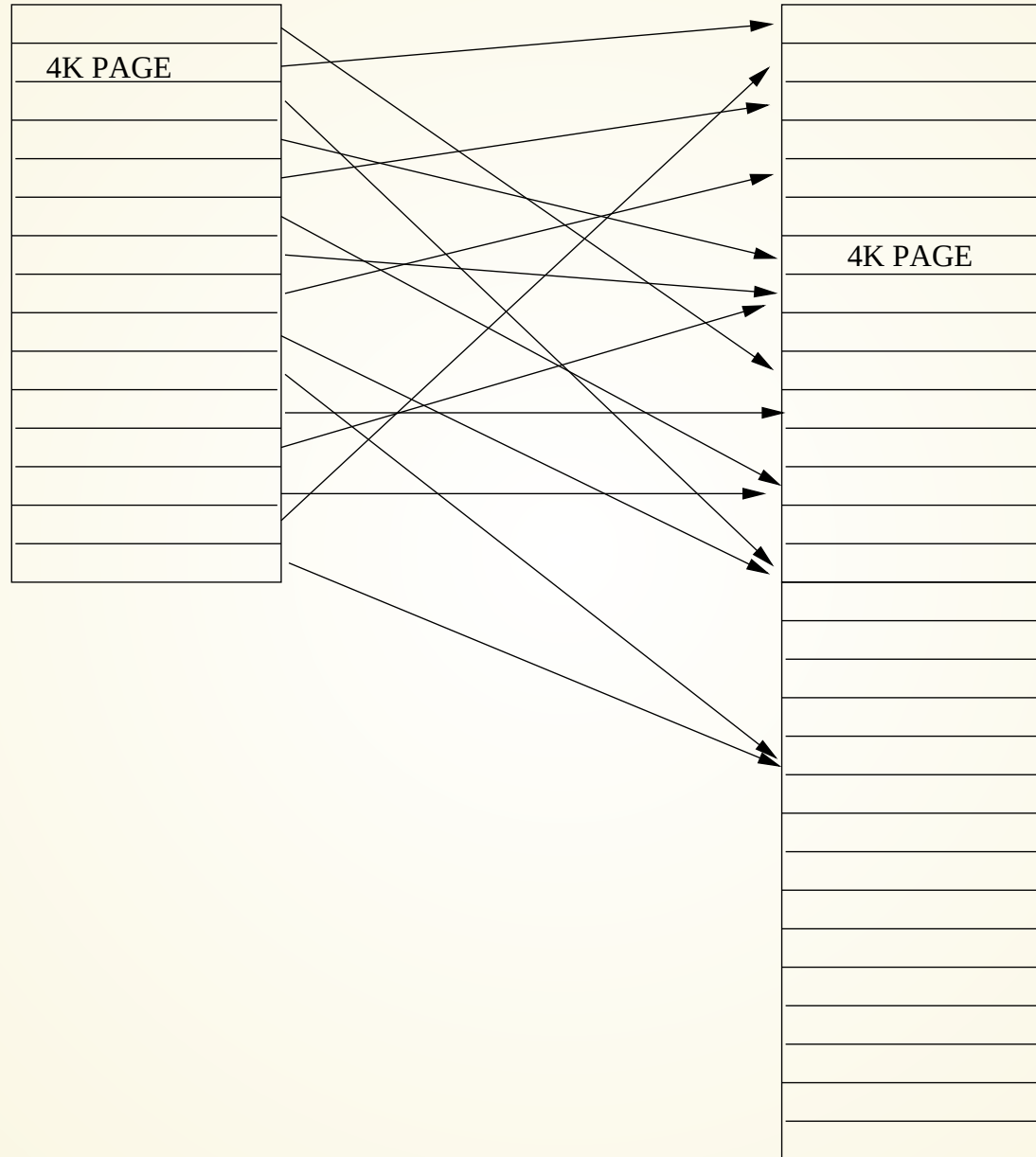


RECAP : MMU



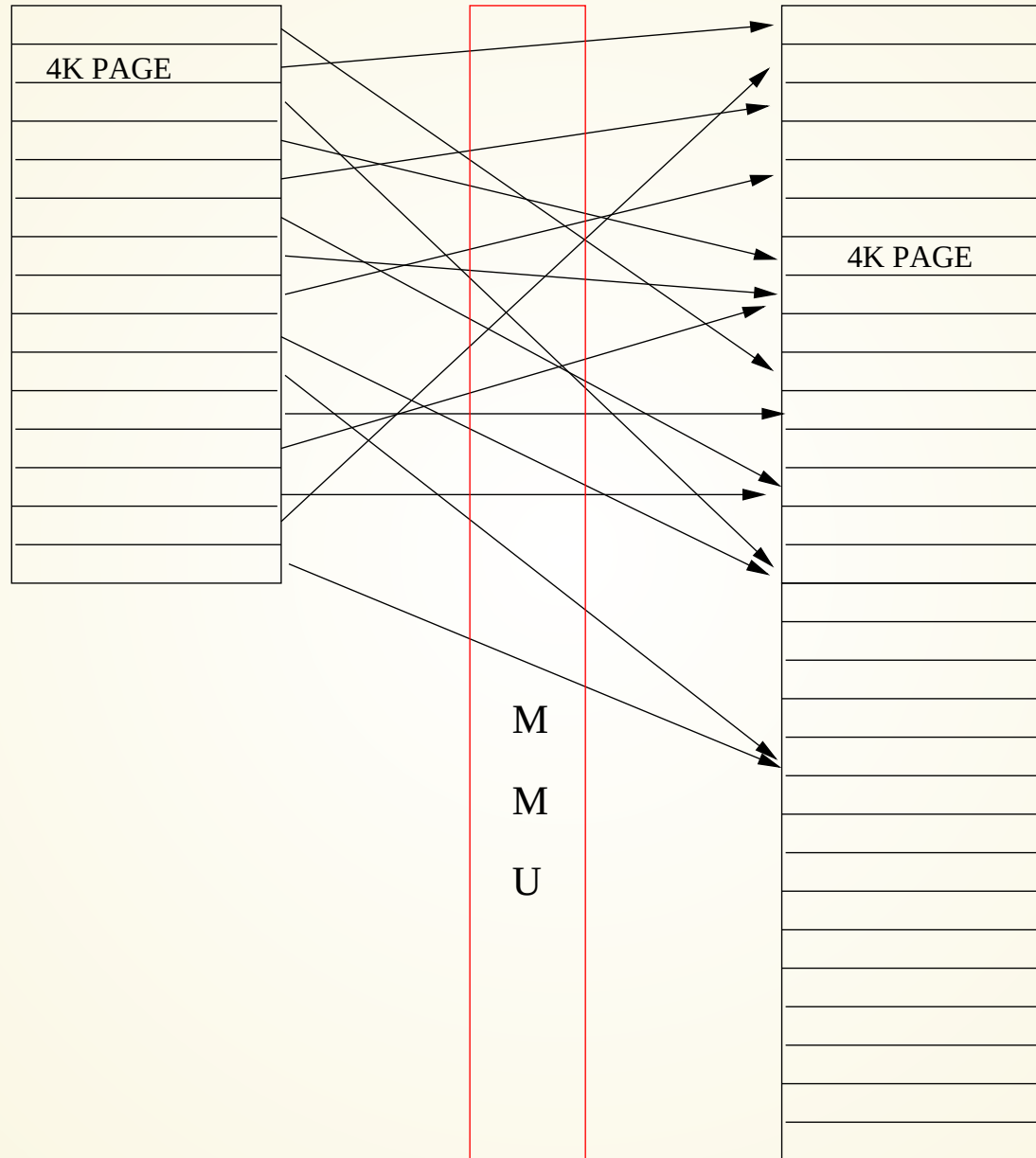
VIRTUAL ADDRESS SPACE

PHYSICAL ADDRESS SPACE



VIRTUAL ADDRESS SPACE

PHYSICAL ADDRESS SPACE

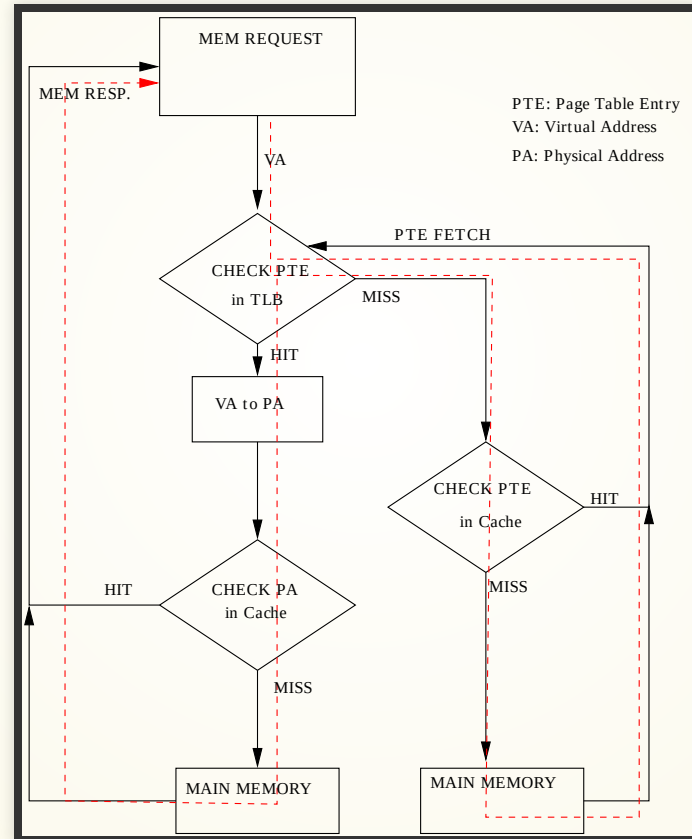


RECAP: MMU OPERATION

- Translation Lookaside Buffer
 - Keeps a page table for virtual to physical address translation.
 - 4GB memory with page size of 4K => ~4MB
 - Each process has a different page table.
 - page table is kept in main memory.
 - Each access will need two accesses to main memory.
 - TLB acts as a cache for page table entries (PTE).

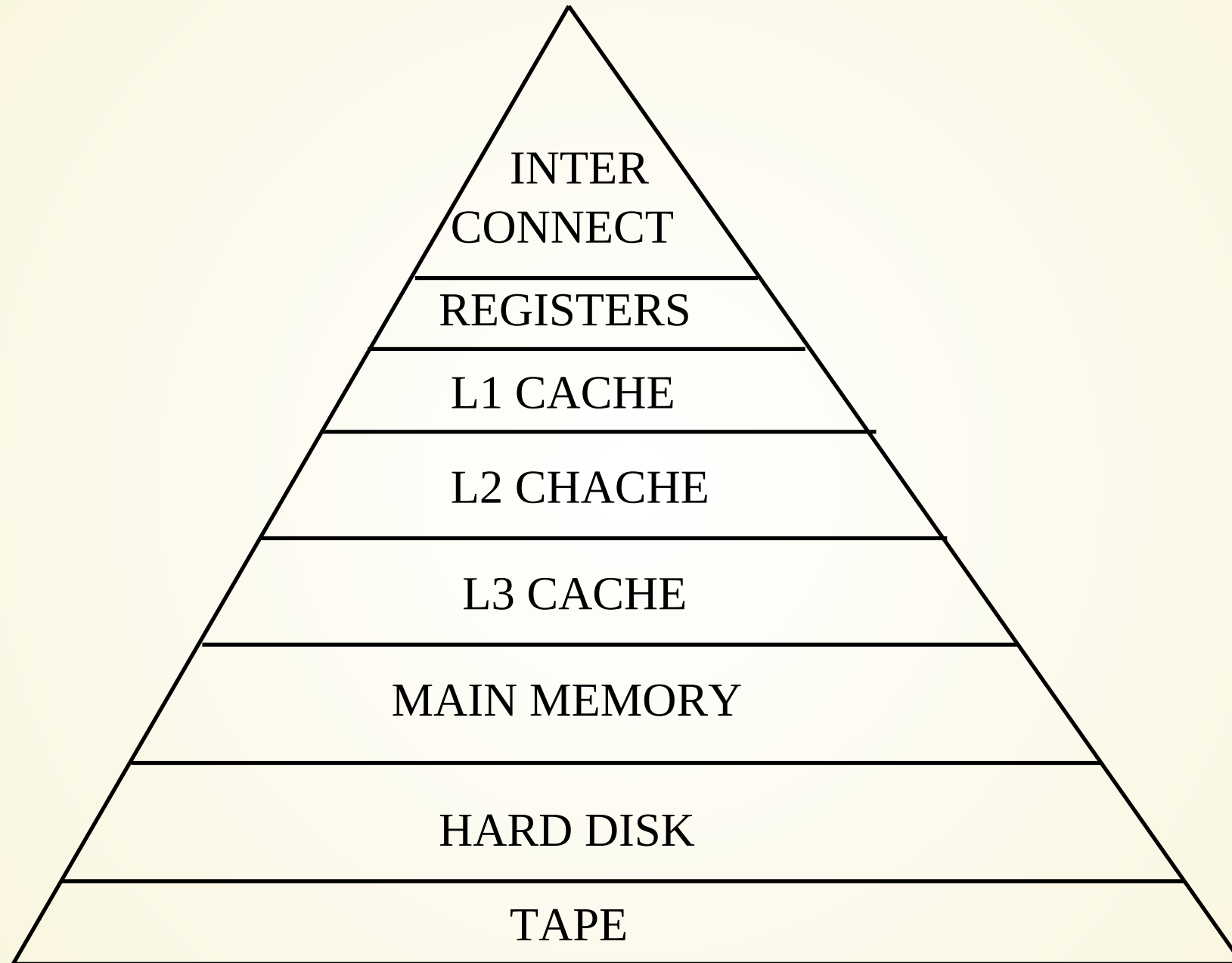


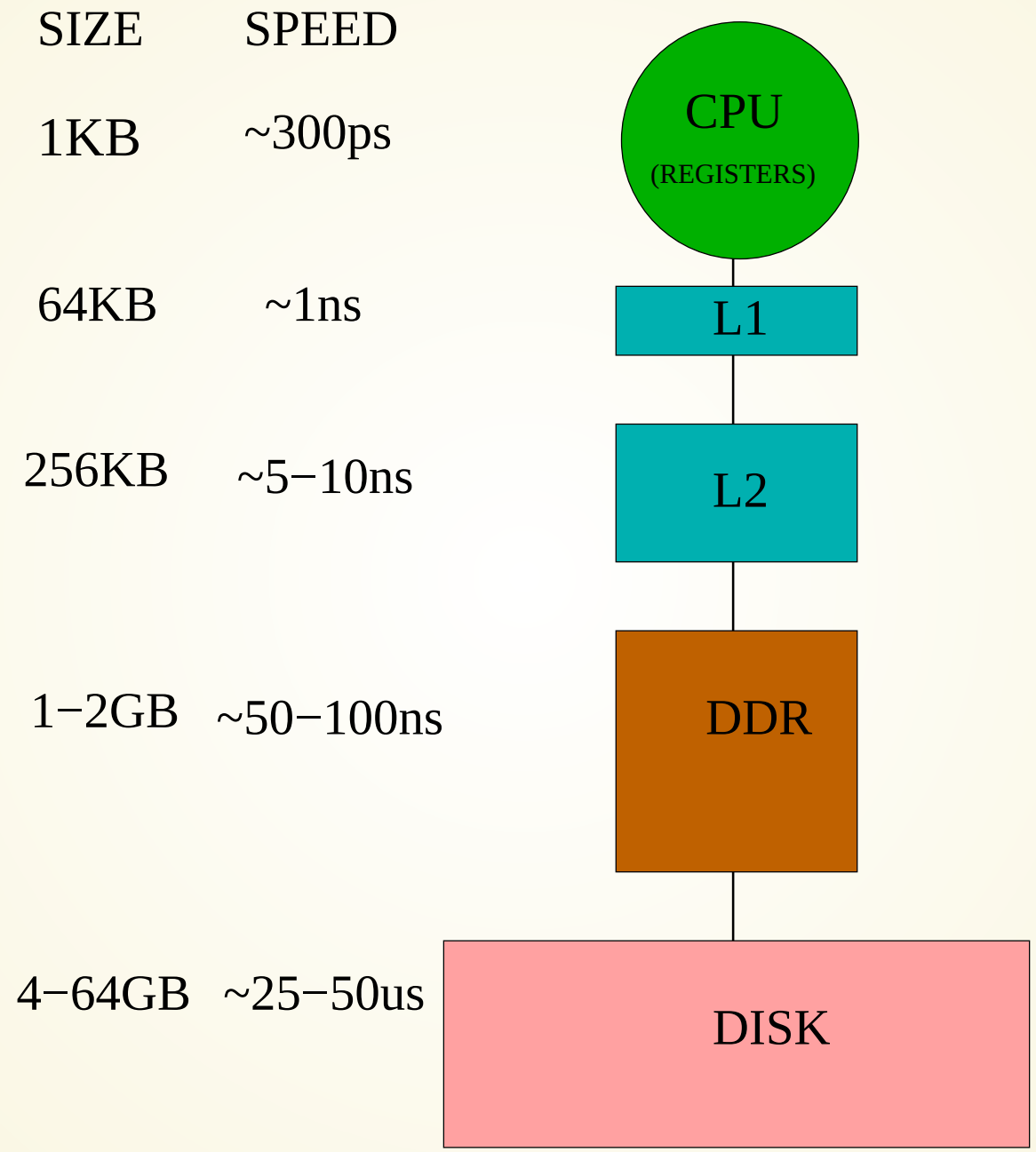
RECAP : LIFE OF A MEMORY REQUEST

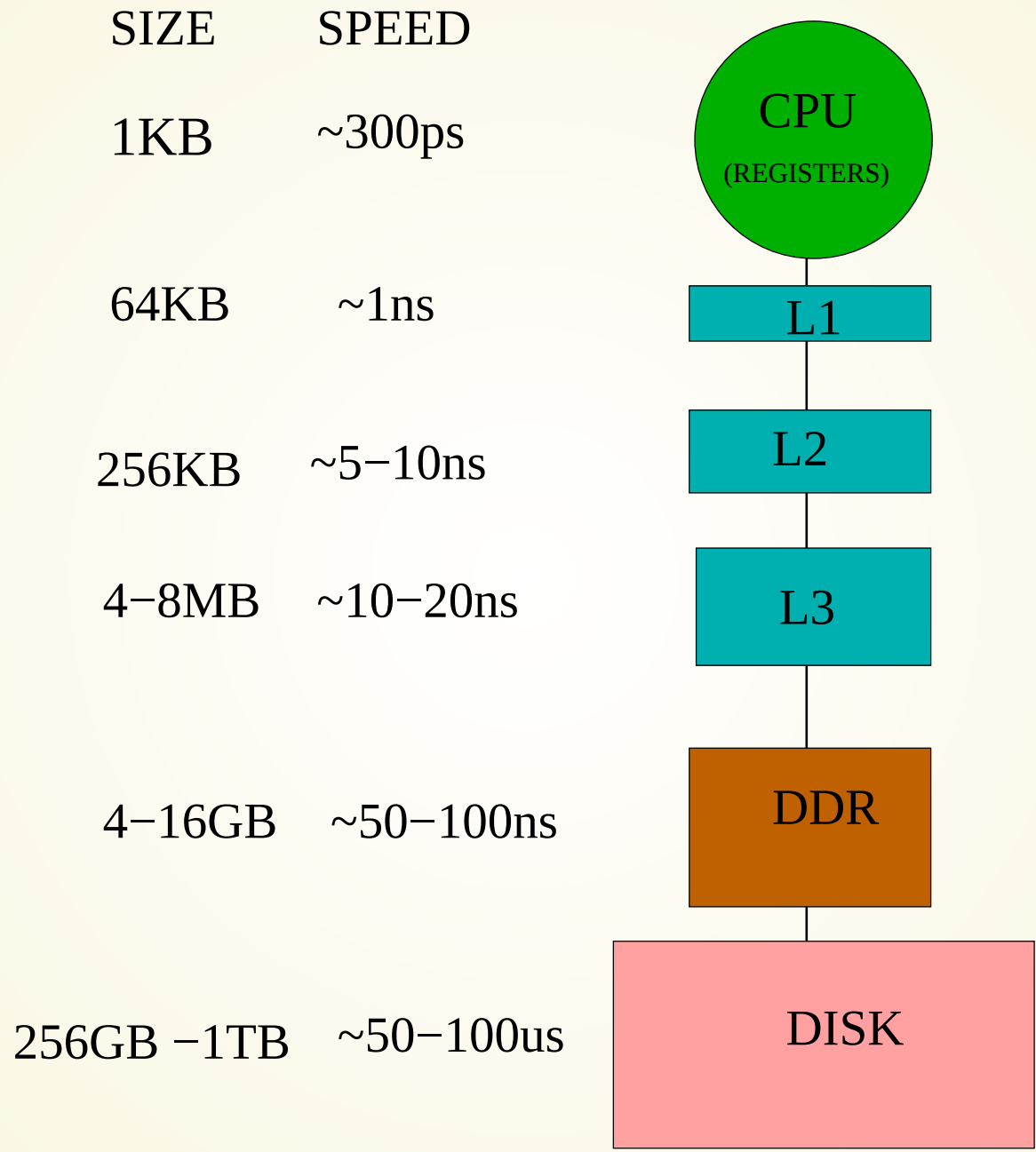


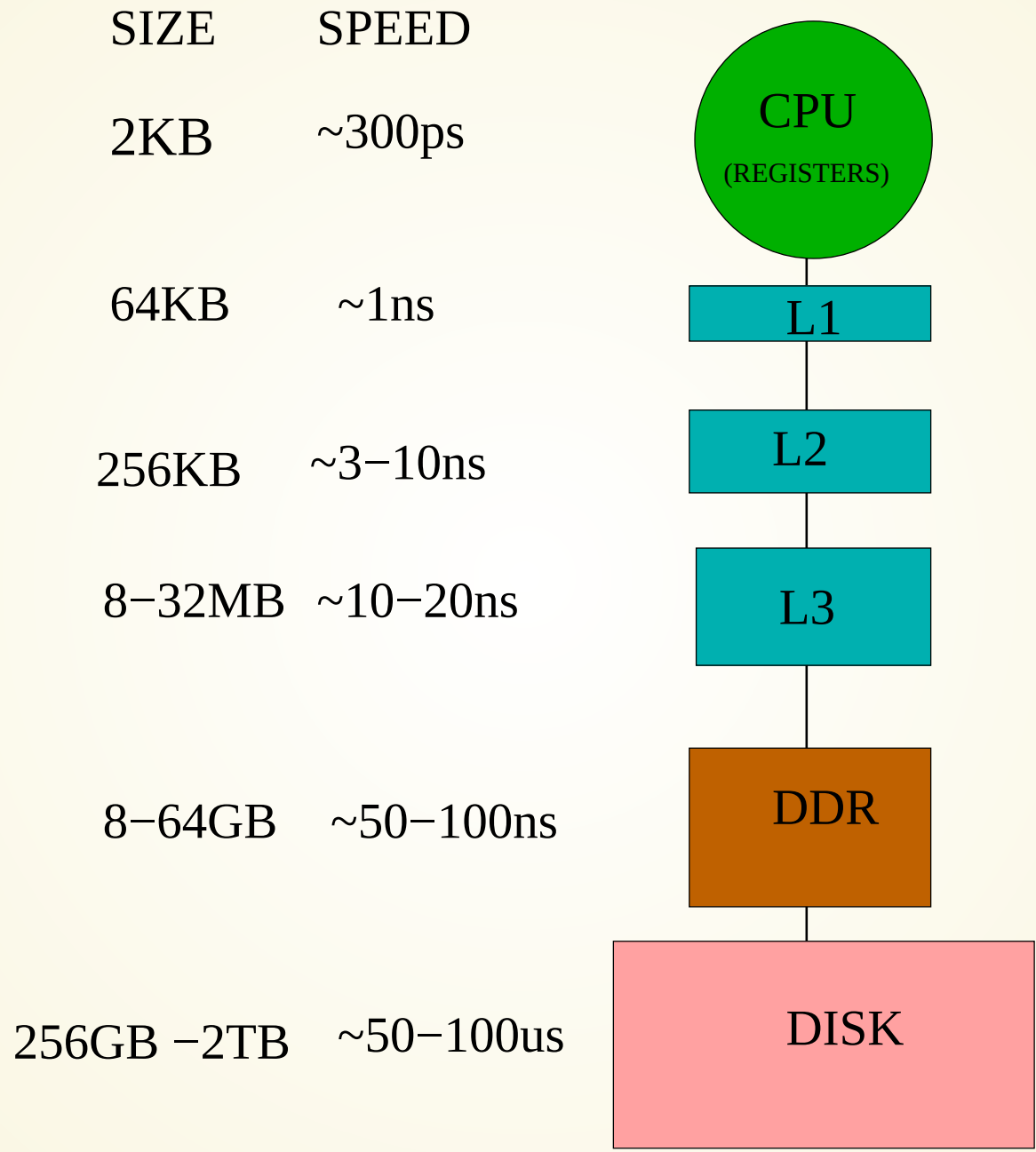
RECAP : MEMORY HIERARCHY



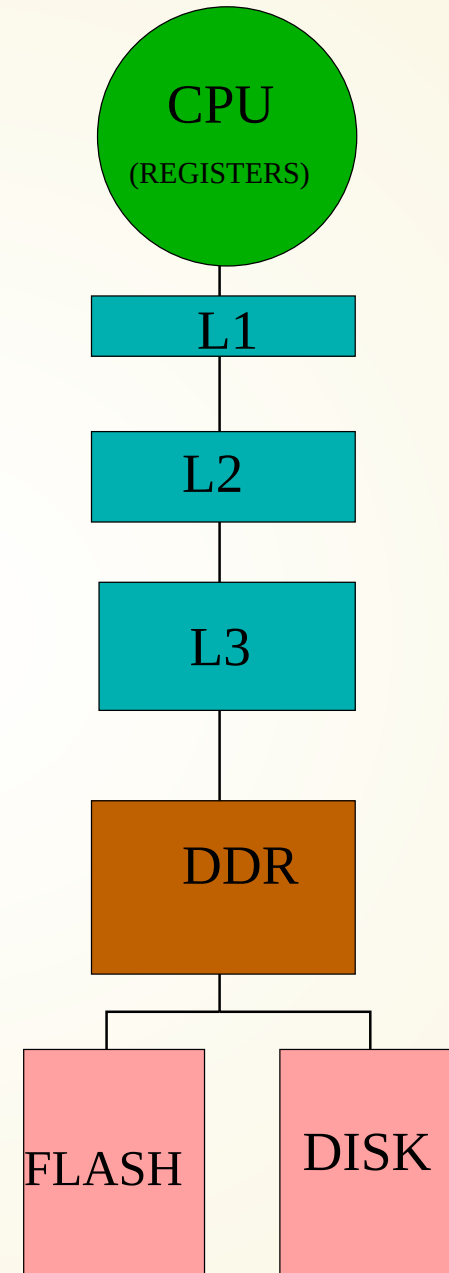








SIZE	SPEED
4KB	~300ps
64KB	~1ns
256KB	~3-10ns
16-64MB	~10-20ns
32-256GB	~50-100ns
16-64TB	~5-10ms
1-16TB	~100-200us



RECAP: PROCESSORS



	Single Data	Multiple Data
Single Instr.	SISD	SIMD
Multiple Instr.	MISD	MIMD



	Single Data	Multiple Data
Single Instr.	CPU e.g 8086	SIMD
Multiple Instr.	MISD	MIMD



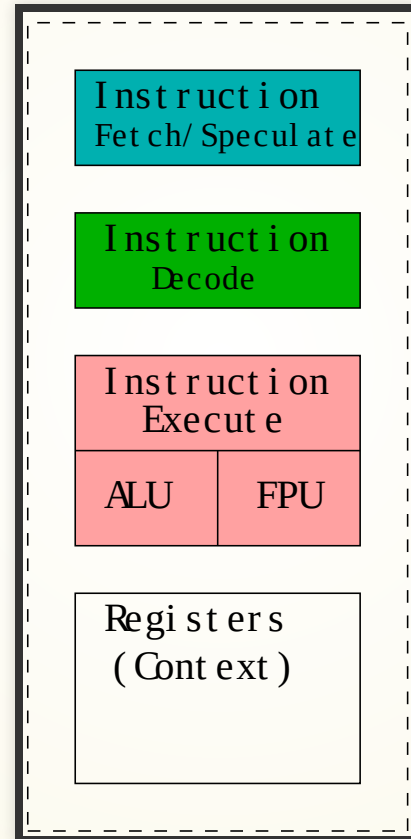
	Single Data	Multiple Data
Single Instr.	CPU e.g 8086	VPU, GPU e.g CRAY NVIDIA
Multiple Instr.	MISD	MIMD



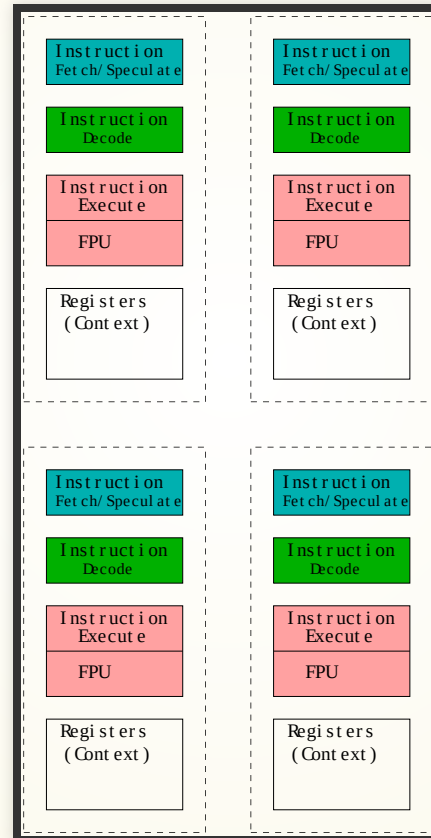
	Single Data	Multiple Data
Single Instr.	CPU e.g 8086	VPU, GPU e.g CRAY NVIDIA
Multiple Instr.	MISD	Multicore e.g Intel i7



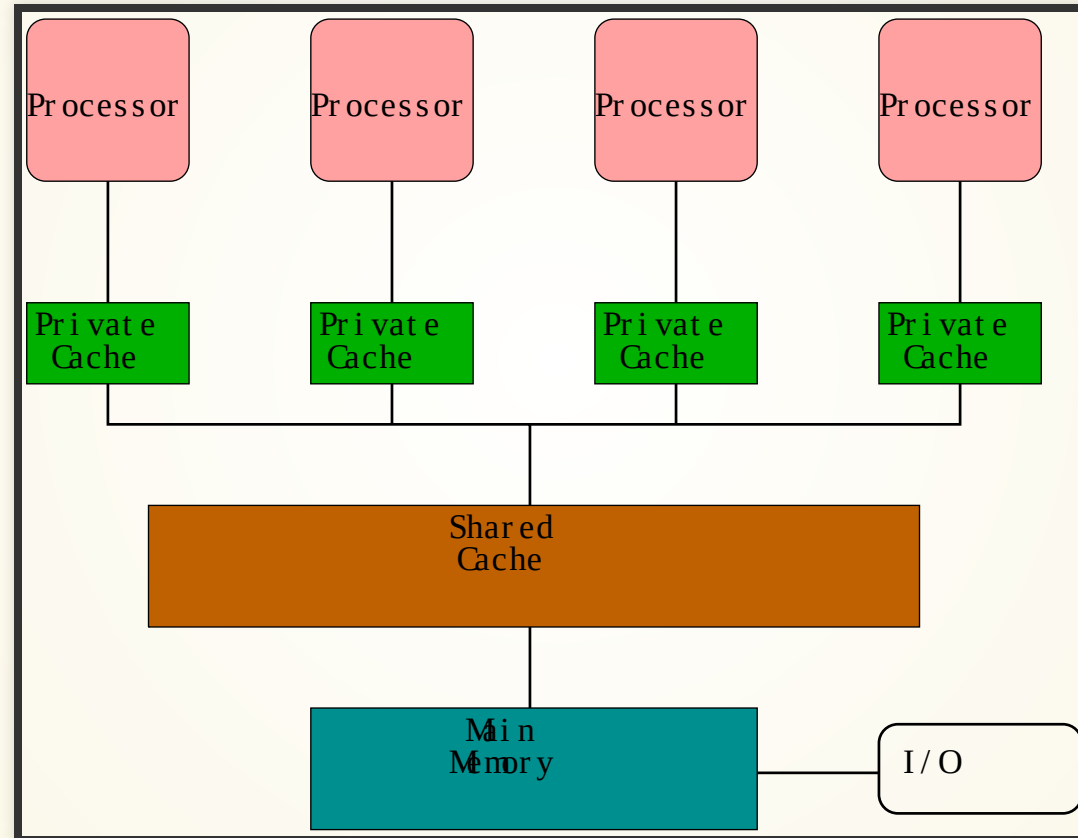
UNIPROCESSOR



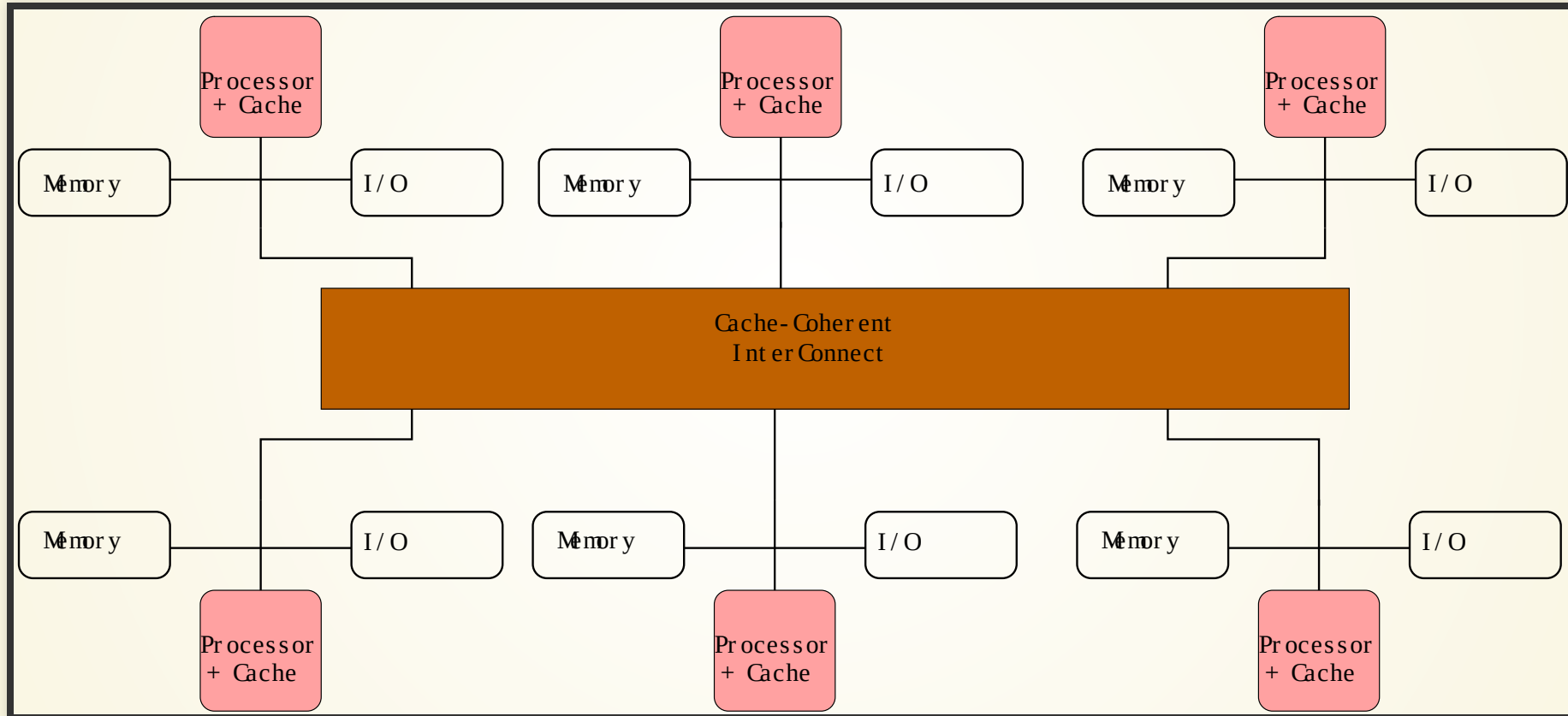
MULTIPROCESSOR



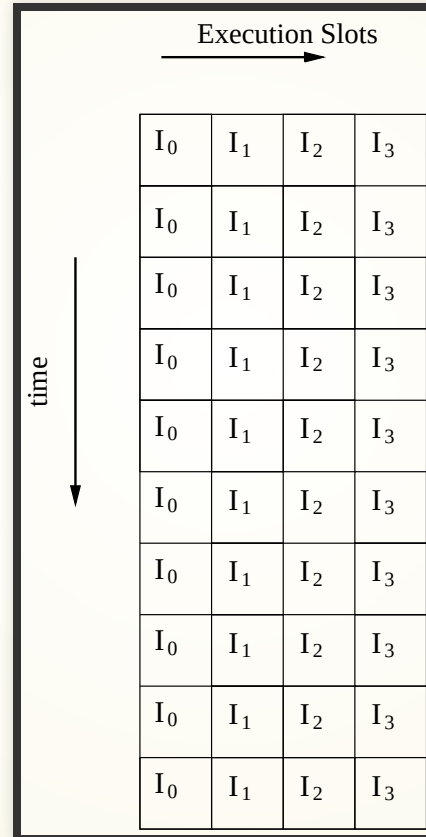
MULTIPROCESSOR: UMA



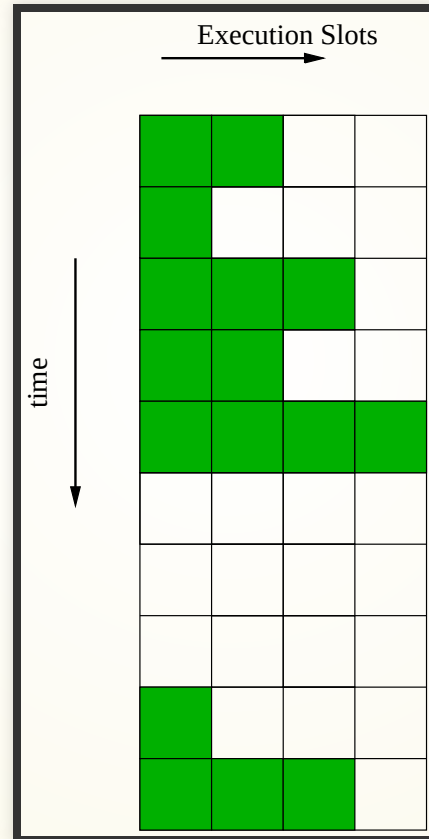
MULTIPROCESSOR: NUMA



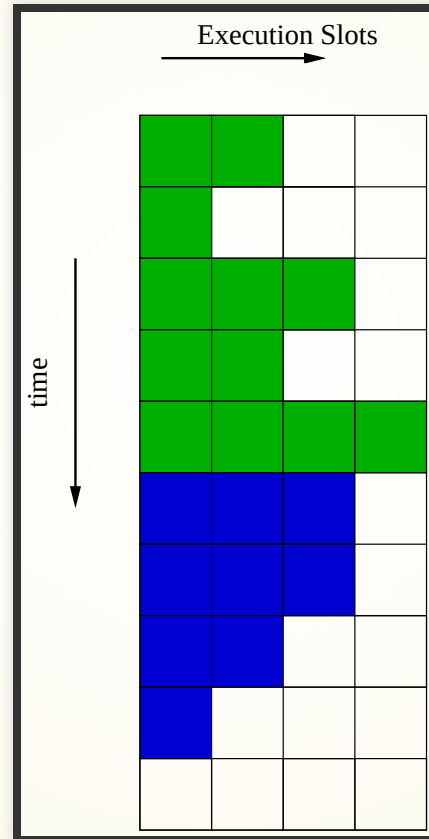
PROCESSOR: VLIW



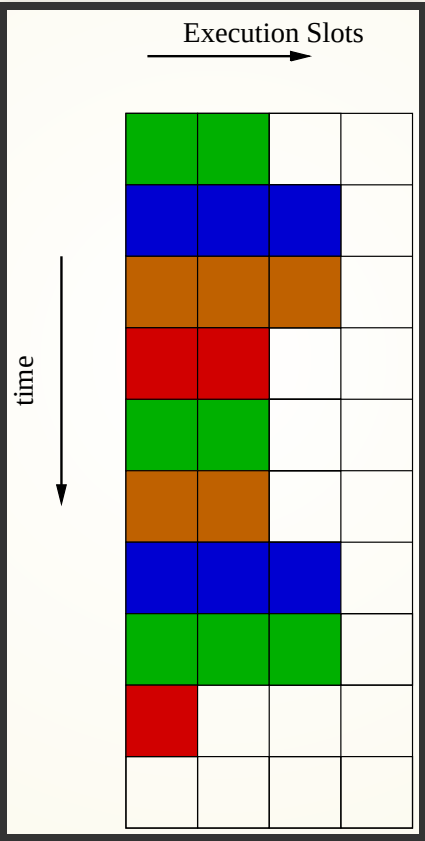
MULTI-THREADING: SUPERSCALAR



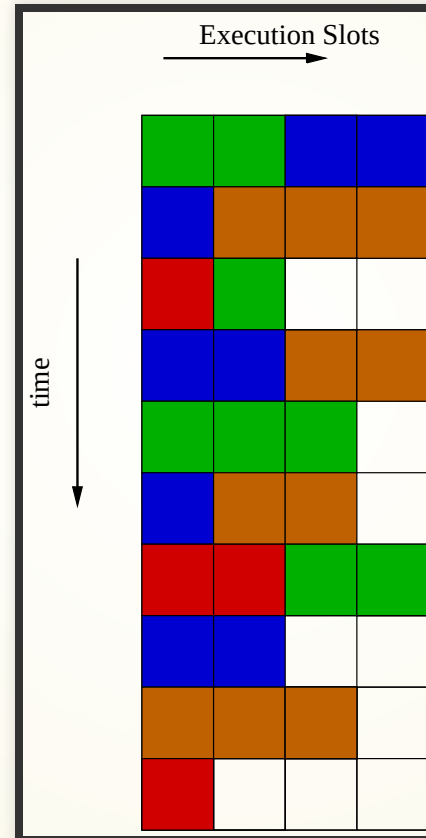
MULTI-THREADING: COARSE-GRAINED



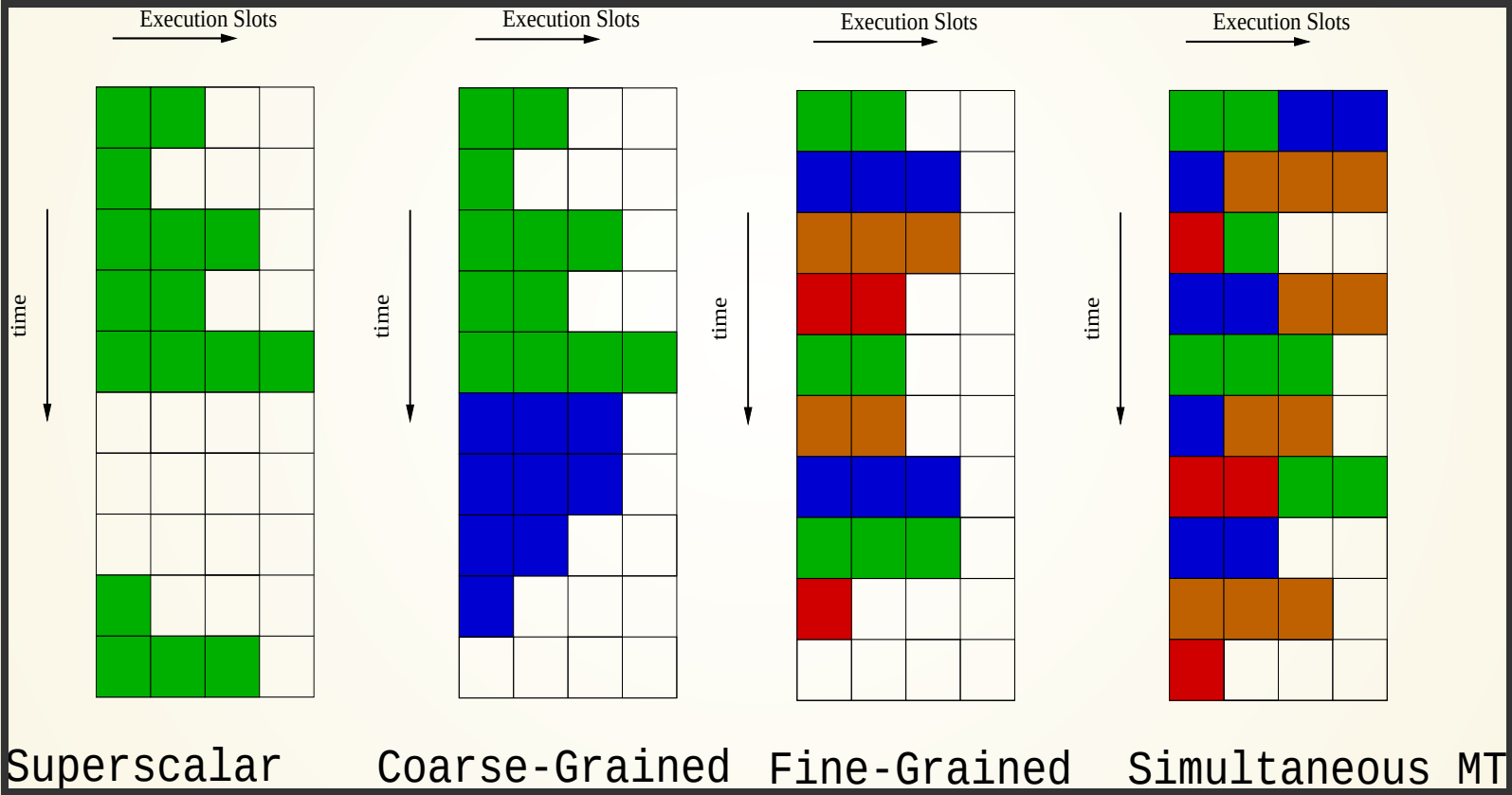
MULTI-THREADING: FINE-GRAINED



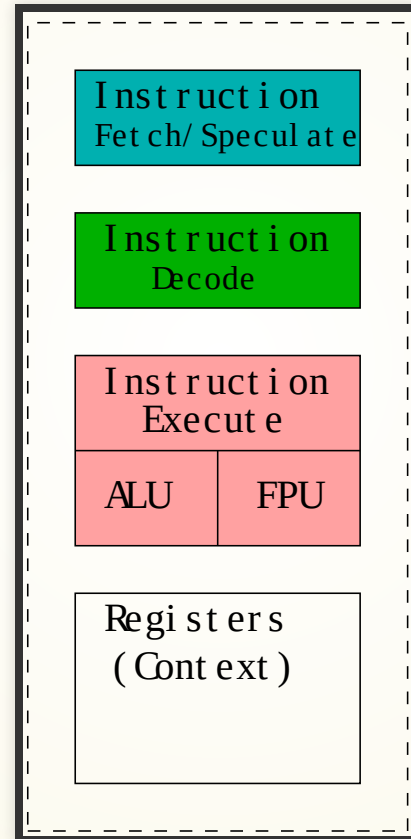
MULTI-THREADING: SIMULTANEOUS



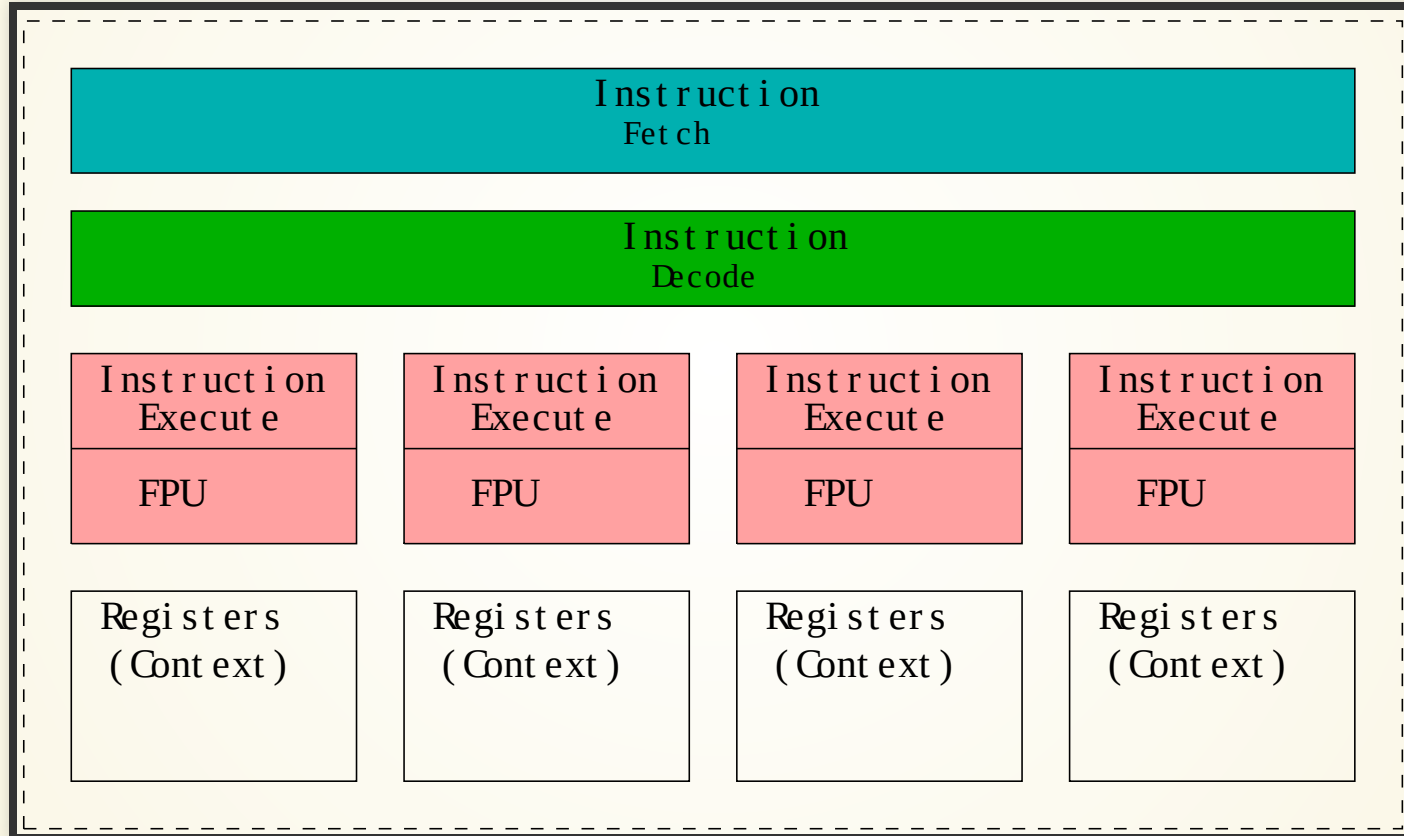
REACP: MULTI-THREADING



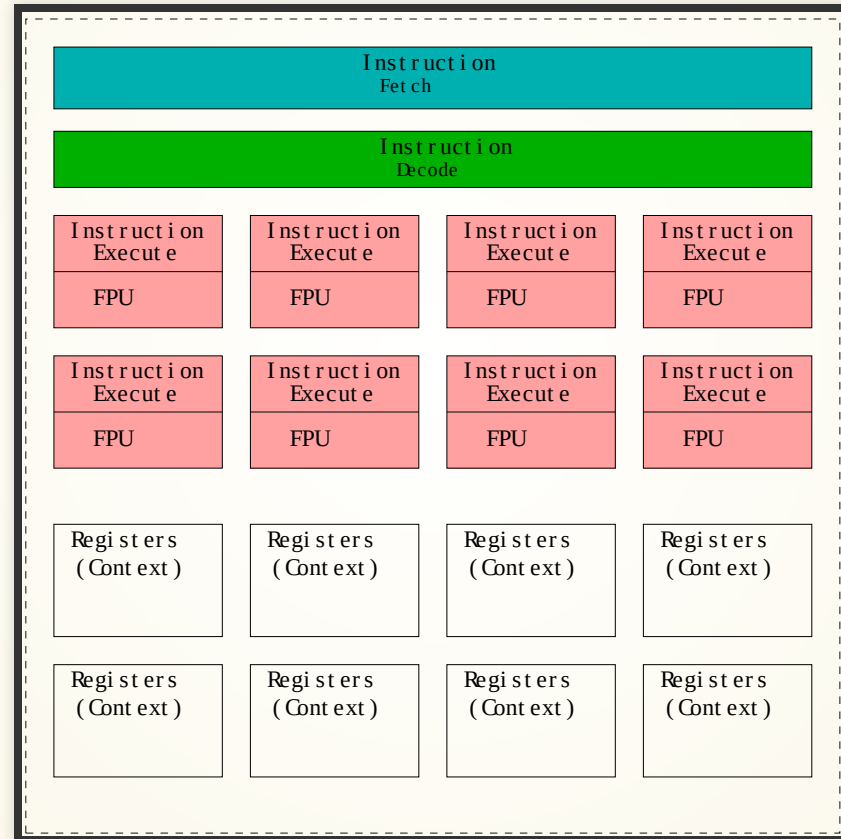
GPU ARCHITECTURE : UNIPROCESSOR



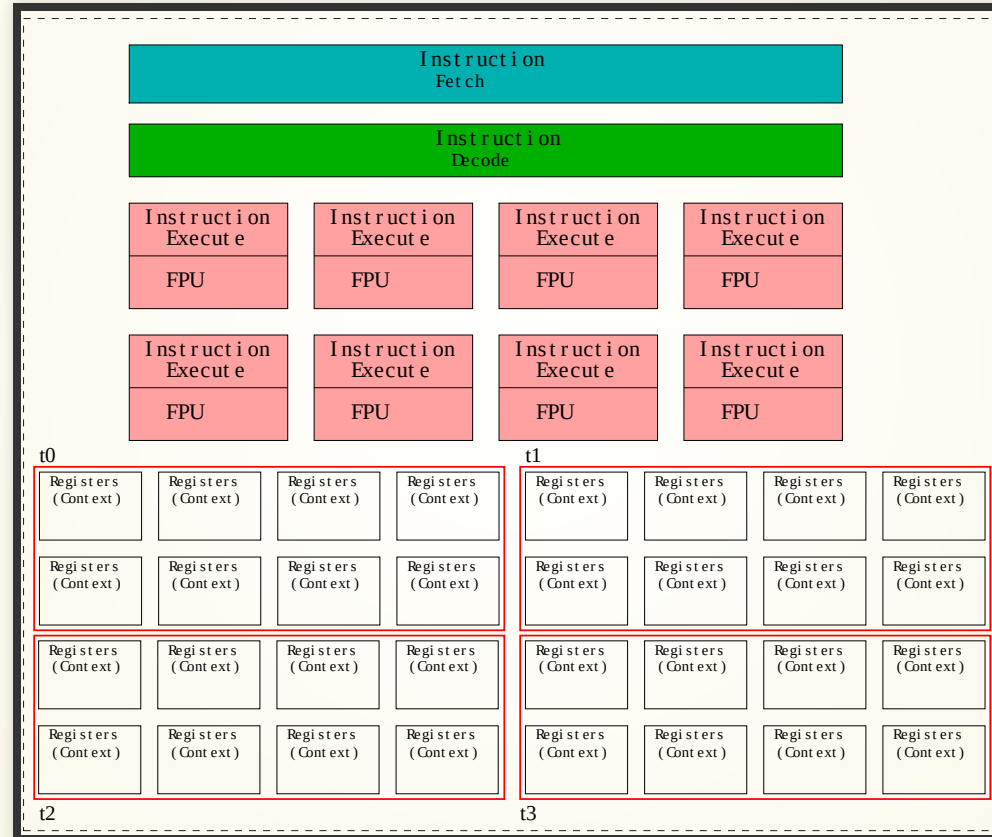
GPU ARCHITECTURE : SIMD



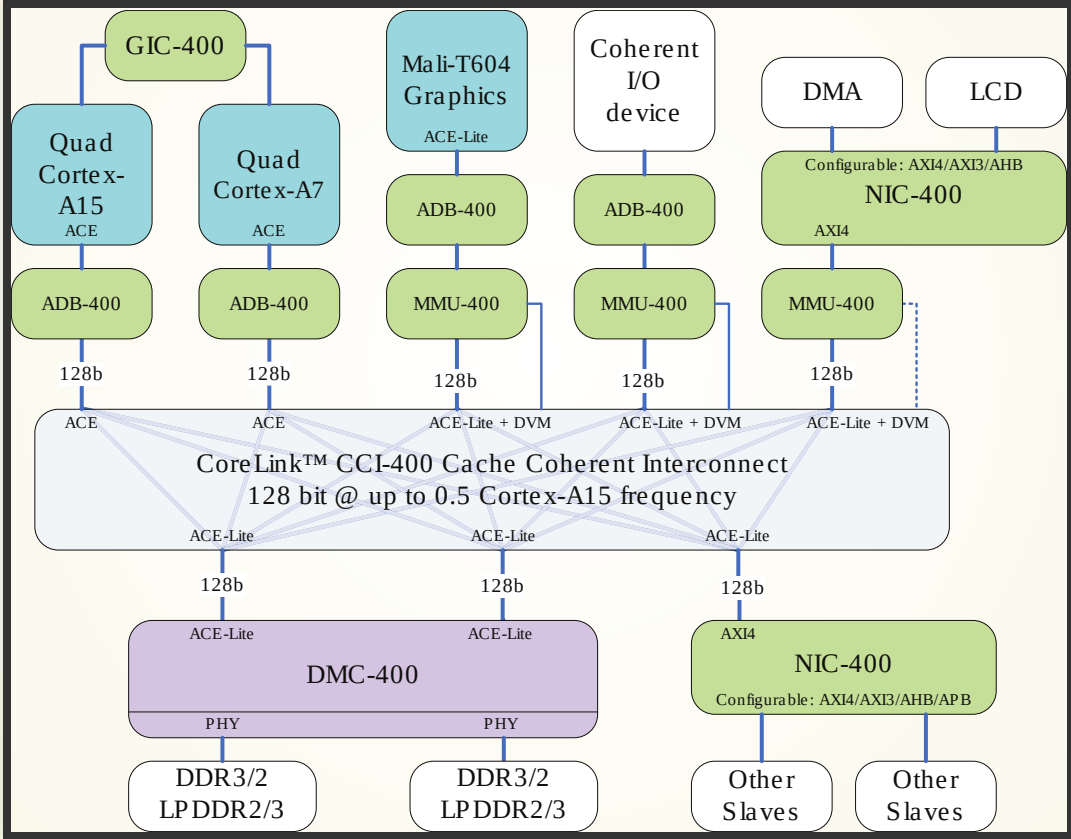
GPU ARCHITECTURE : SIMD



GPU ARCHITECTURE : SIMD



EXAMPLE HETEROGENEOUS SOCS



DOMAIN SPECIFIC ARCHITECTURE

