

Lecture 17: Cache performance and virtual memory

Tuesday, March 6, 2018 10:15 AM

Outline

- Cache performance
 - AMAT
- Multi-level caches
- Real examples
- Virtual memory
 - Why?
 - Multi-level page tables
- Improving VM performance

System performance with caching

of misses

Average memory access time

Miss rate/ratio = $\frac{\# \text{ misses}}{\# \text{ accesses}}$

Misses per instruction

50 Mpi → .05 misses/instr

Time for application to run

CPI

Time pr app = (CPU exec. cycles + memory stall cycles) · cycle time

↳ assumes blocking cache

↳ $\frac{\# \text{ memory accesses}}{\text{program}} \cdot \text{miss rate} \cdot \text{miss penalty}$

how to reduce miss rate? →

- ↳ incr. cache capacity
- ↳ use diff replacement policy
- ↳ incr associativity

Row major

A spatial locality

B

temporal locality

to incr. locality work on a block

Better temporal locality

Pick "right" block size

B^T ← column

Average memory access time (AMAT)

AMAT = time for a hit + Miss ratio · miss penalty

= 2 + 0.01 · 100 = 3 cycles

(1 - miss) · 2 + miss · (2 + 100)

Double freq

miss penalty → 200 cycles

AMAT = 2 + 0.01 200 = 4 cycles

→ hit time 2 cycles

→ miss ratio of 0.01

miss penalty of 100 cycles

Hit	Miss
lookup tags	lookup cache
2 cycles	2 cycles
	send mem pg
	100 cycles

Multi-level caches

$$AMAT = \text{hit} + \text{miss ratio} \cdot \text{miss penalty}$$

$$= \text{hit}_{L1} + \text{miss ratio}_{L1} \cdot (\text{hit}_{L2} + \text{miss ratio}_{L2} \cdot \text{miss penalty})$$

$$= 2 + 0.01 \left(\frac{10}{10} + \frac{0.05 \cdot 200}{20} \right)$$

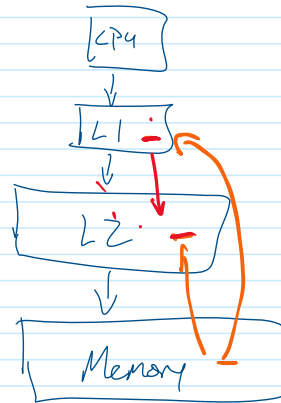
$$= 2 + 0.2$$

$$= 2.2 \text{ cycles}$$

$AMAT_{L2}$

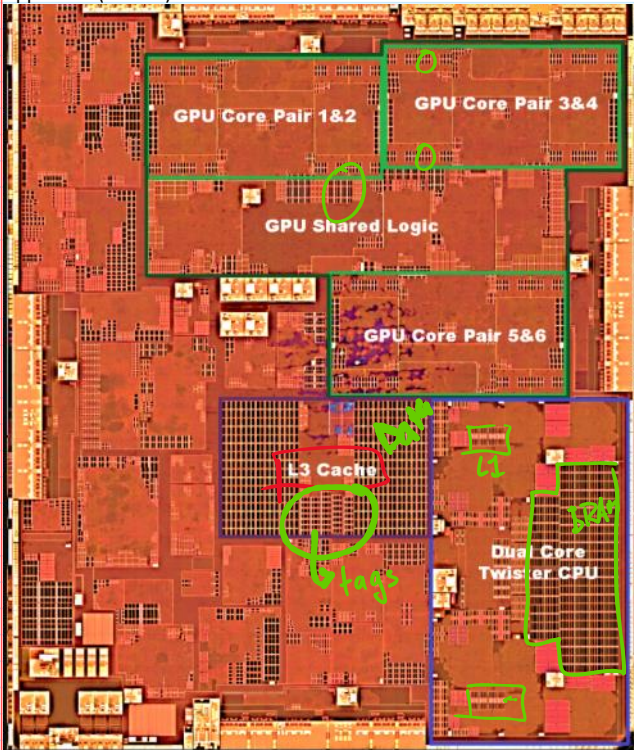
to reduce \rightarrow add a cache

$L2$ miss ratio of 5%
hit time of 10 cycles

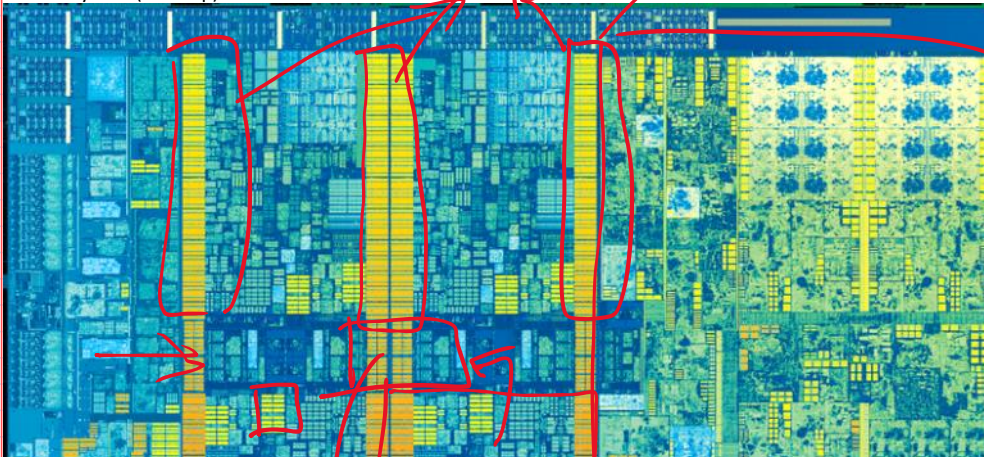


Cache examples

Apple A10 (mobile)

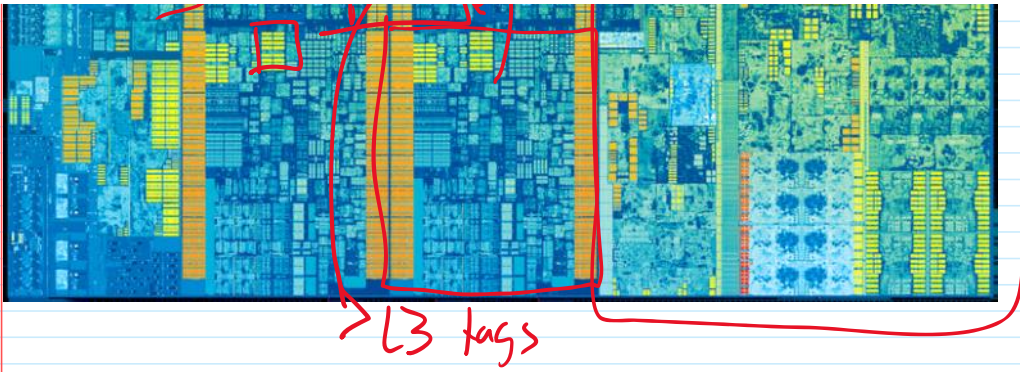


Intel Kaby lake (desktop)

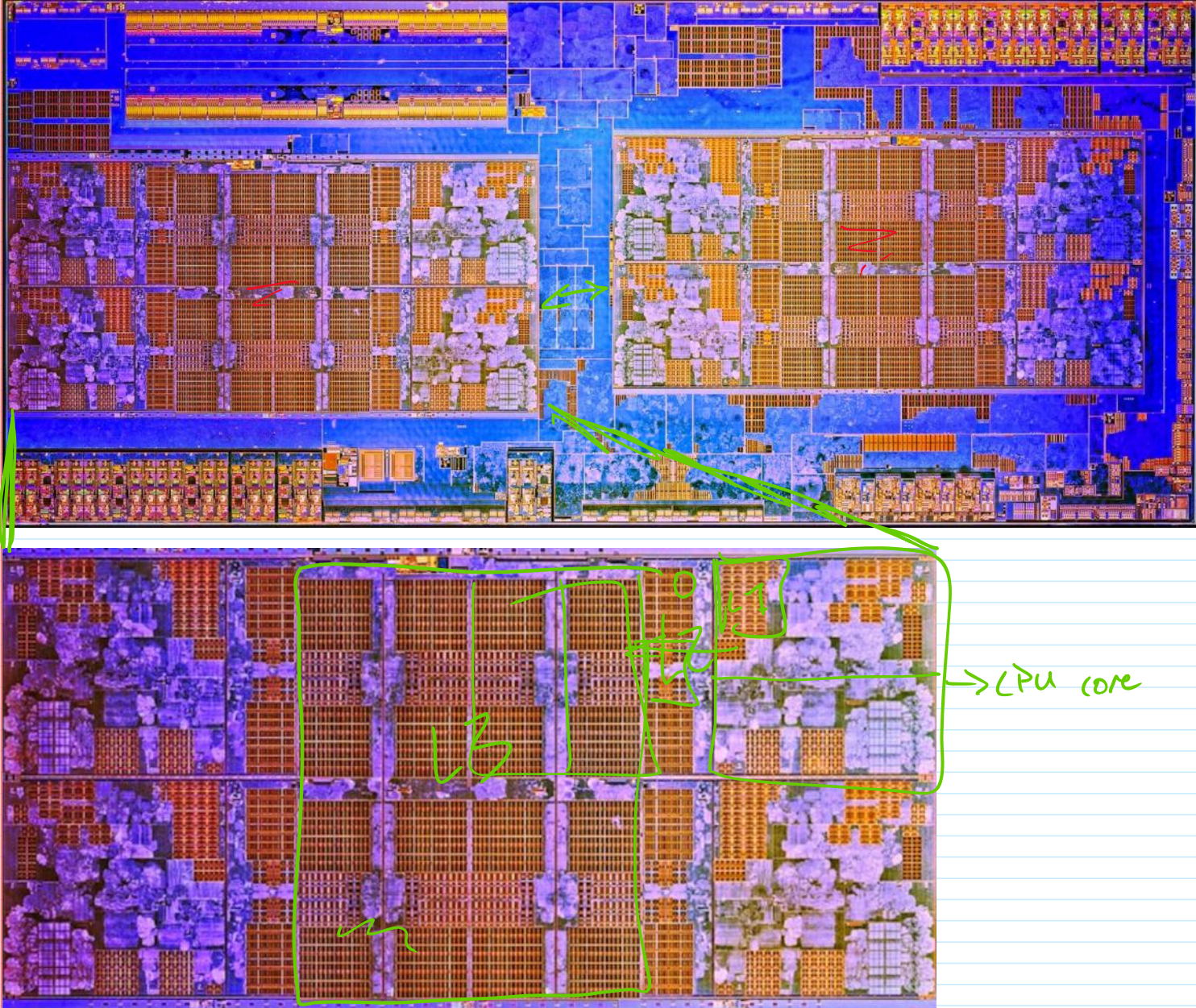


8MB $L2 \rightarrow 256K$

GPU



AMD Epyc (Server)



Virtual Memory

What is VM?

Why?

- illusion to programmer that she has more memory
- security & flexibility in assigning physical memory
- combat fragmentation
- programs appear to be running on their own machine
- part of ISA o/w - b/w parsize
 - ↳ 32 bit ISA
 - ↳ 64 bit
- originally looked like more memory by being a cache of hard drive
 - ↳ swapping

Like a cache

Block size? → 4KB or 2MB

↳ page size

or physical frame

associativity? fully associative b/c → any page can be mapped anywhere in virtual address space

Capacity? size of main/physical memory

