



CS107 Lecture 25

Optimization, Caching, Writing Cache-Friendly Code

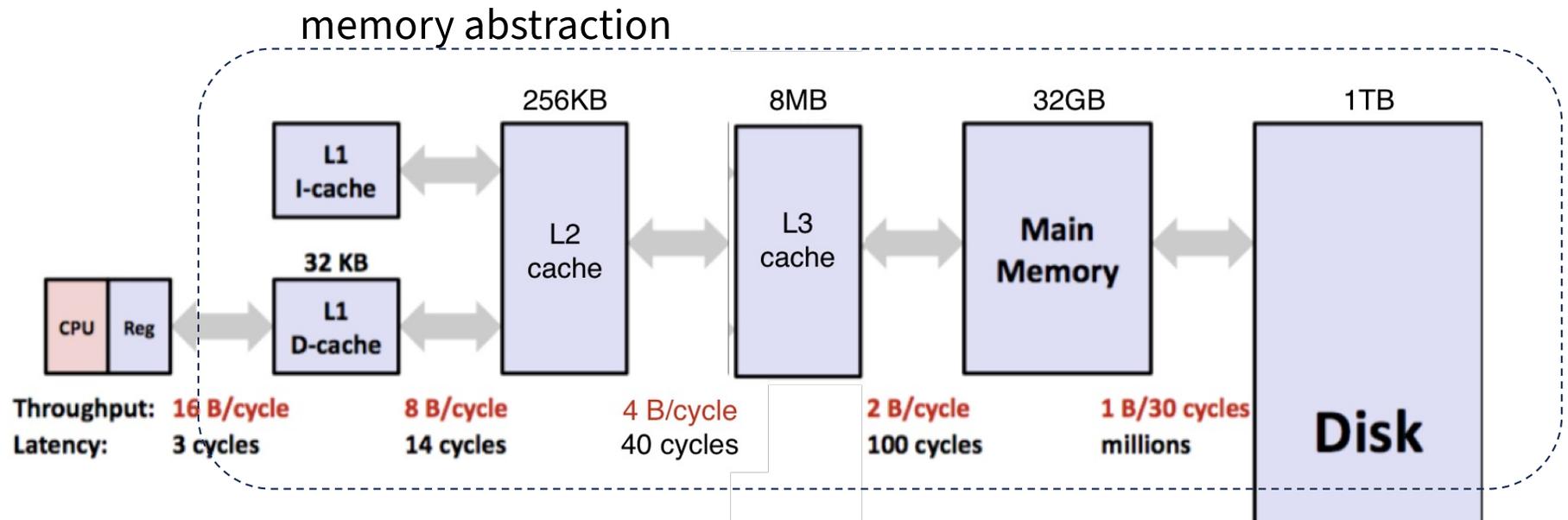
Reading: B&O 5

Caching

Processor speed is **not the only bottleneck** in program performance—
memory access is perhaps even more of one!

Memory exists in levels and goes from **very fast** (registers) to **very slow** (disk).

Frequently accessed data tends to stay in faster memory.



Caching

Virtually all caching depends on **locality**.

Temporal locality

- Repeat access to the same data tends to be co-located in **time**
- **Intuition**: Memory I've accessed recently is likely to be accessed again very soon.

Spatial locality

- Related data tends to be collocated in **space**.
- **Intuition**: Memory I've recently accessed might be right next to memory I'll need very soon.

Realistic scenario:

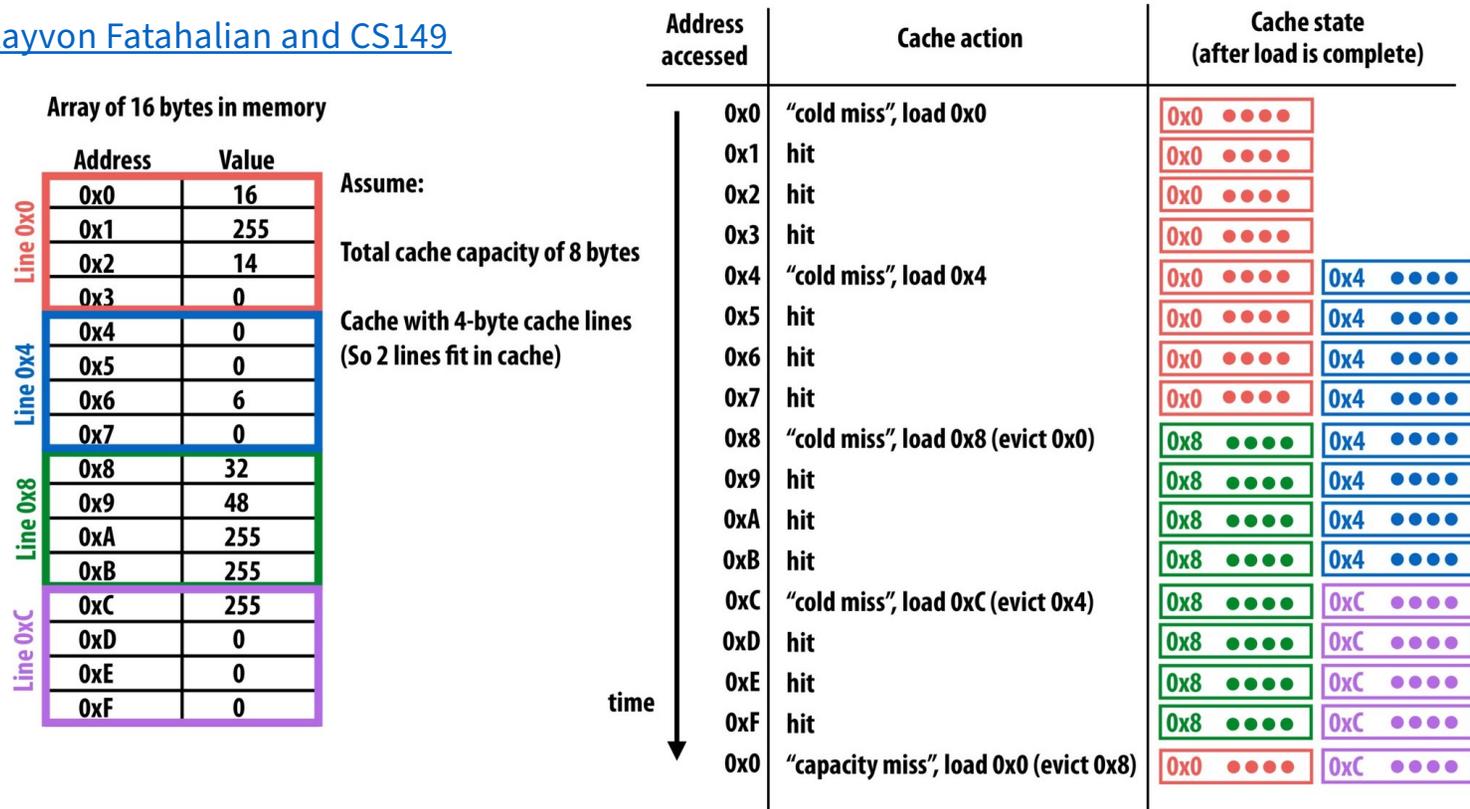
97% cache hit rate, cache hits incur 1 clock cycle, whereas cache misses incur 100

Thought Question: How much time is spent on the 3% of accesses that are cache misses?

Caching

Virtually all caching depends on **locality**.

Credit: [Kayvon Fatahalian and CS149](#)



Caching

The CPU can **execute instructions quickly**, but if it must wait for memory, it sits idle. Caches keep data the CPU needs **as close to it as possible**. Note: RAM officially holds the **data** and caches store **copies**.

Demo: cache.c

