

Lecture 14: Caching intro

Thursday, February 22, 2018 9:32 AM

Outline

- Why memory hierarchy
- How to control memory hierarchy
 - Data movement technologies
- Cache basics
 - Direct-mapped
 - Tag & data
 - Miss & hit
 - Handling writes

From last time...

	Typical access time	Cost per gigabyte	Max in machine
SRAM	$\sim 1\text{ns}$	\$300	$\sim 1\text{GB}$
DRAM	$\sim 20\text{ns}$	\$20	$\sim 1\text{TB}$
Flash	$\sim 10\mu\text{s}$	\$0.50	$\sim 100\text{TB}$
Disk	$\sim 10\text{ms}$	\$0.03	$\sim 100\text{TB}$
3rd point	$\sim 100\text{ns} - 1\mu\text{s}$	\$1.25	$\sim 1\text{TB}$

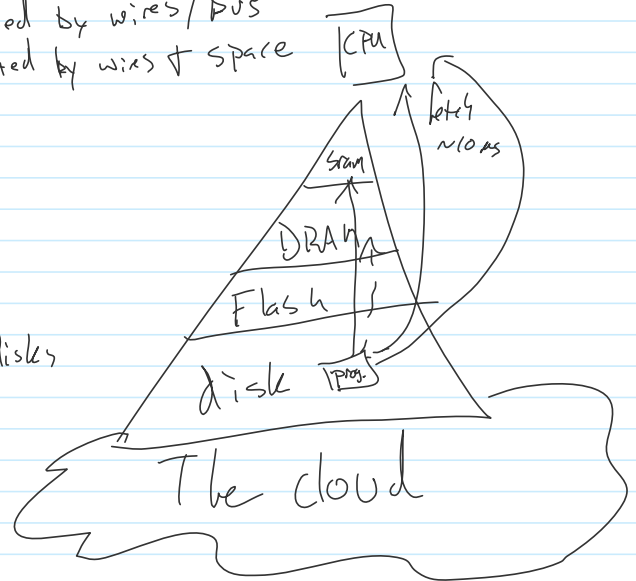
→ limited by power
 → limited by power & space
 → limited by wires / bus
 → limited by wires & space

10x
100x
1000x

not best at anything

Hierarchy

Small amount of SRAM
 Some DRAM
 lots of storage Flash / disks



how to move data?

- caching

→ copy "some"
 ↳ small amount
 grab stuff

temporal locality

↳ reuse same data

near

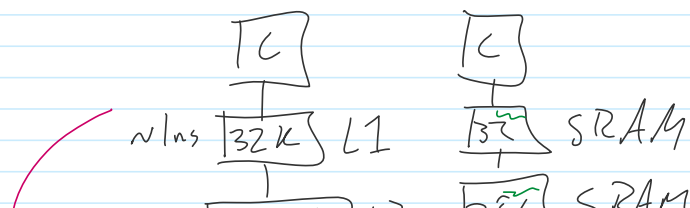
→ things nearby are likely to be used

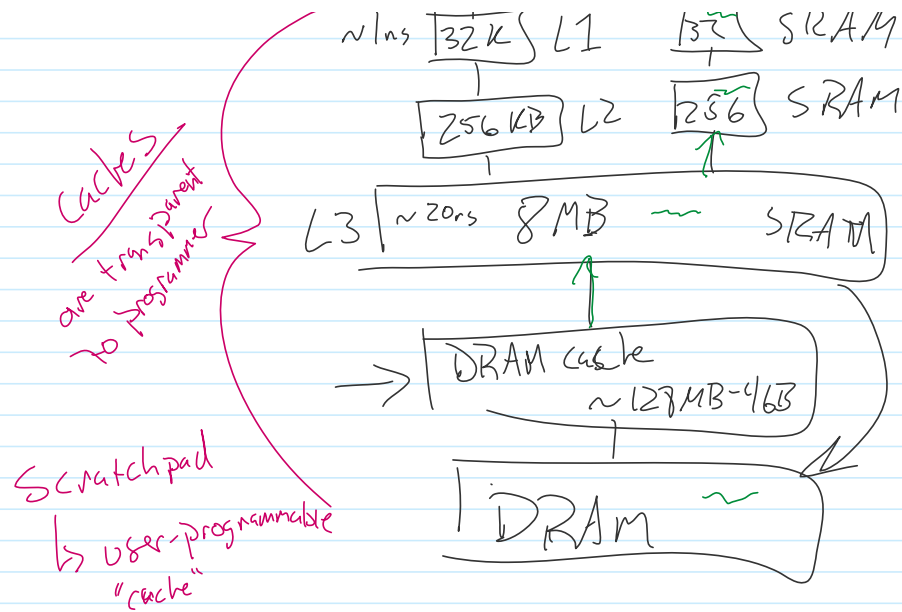
spatial locality

↳ one byte?

move from disk to
 DRAM
 ↳ 64kB 1GB
 ↳ 2MB

move from DRAM
 to SRAM
 ↳ 64 Bytes





Caching

Questions to answer

- Where to look for data?
- How to tell if it's the right data?
- What to do if the data isn't found?
- What to do with writes?
- How to pick a location if multiple options?

take advantage of **locality**

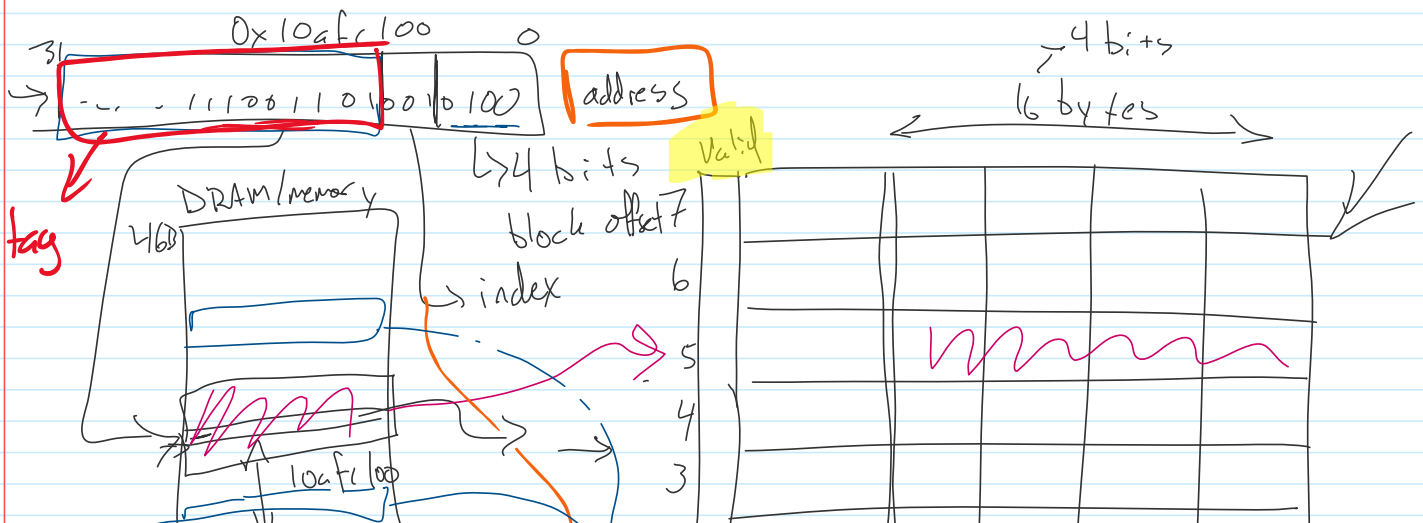
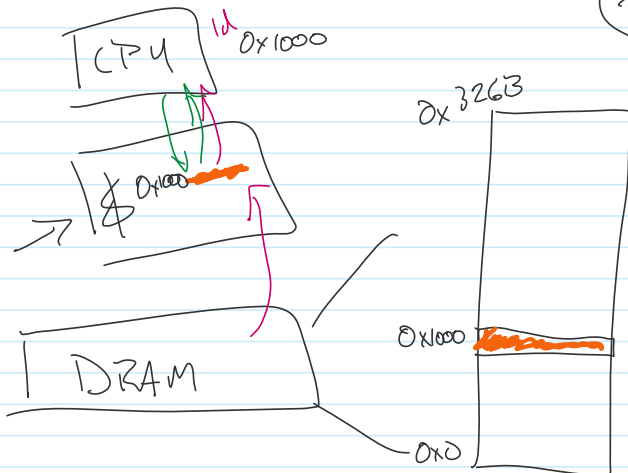
Temporal → things are reused

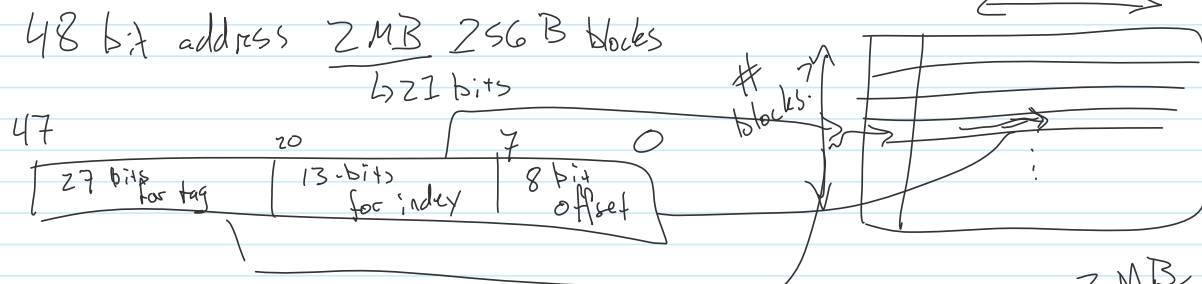
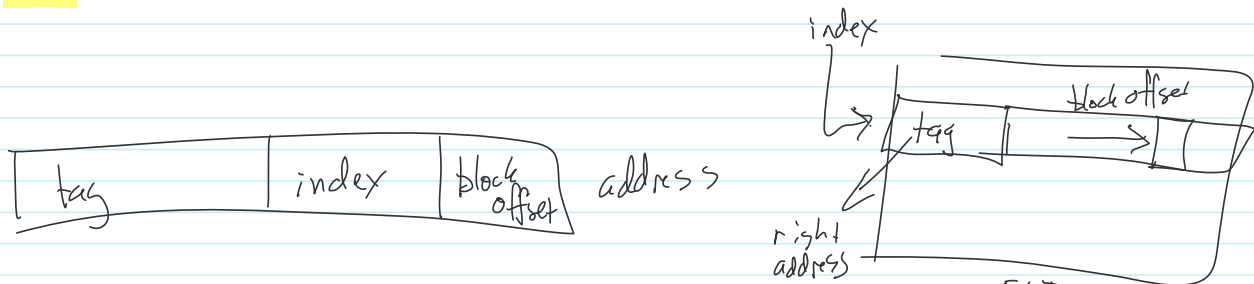
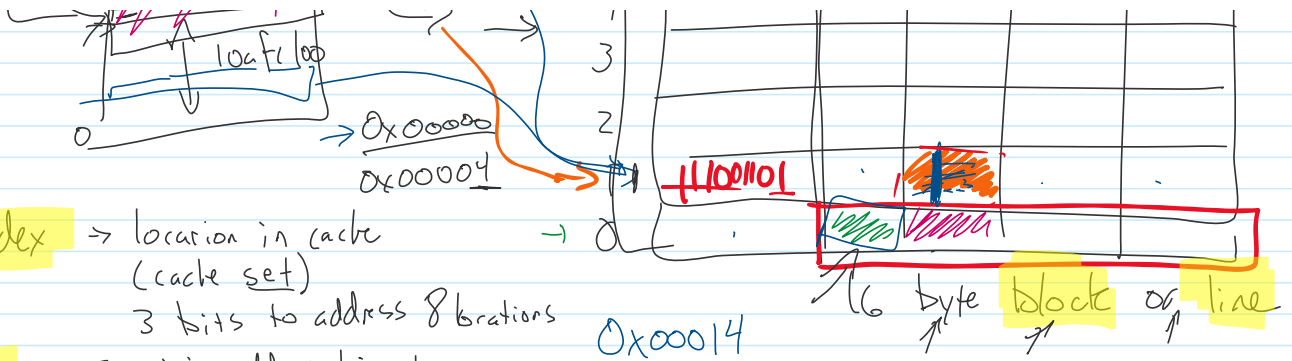
Spatial → likely to access nearby things

keep things around

fetch nearby things

Cache → behave like hash map python dict.





which bits of addr. use for block off., index, tag?

$$\frac{2MB}{256B} \rightarrow \frac{2^{21}}{2^8} = 2^{13}$$