

Lecture 13: Memory technology

Monday, February 19, 2018 1:09 PM

Outline

- Memory technologies
- Memory hierarchy
 - Temporal and spatial locality
 - Blocks/lines
- Packaging technology

Perf. potential of modern CPUs

peak IPC $\sim 4-6$ inst.

freq ~ 4 GHz

per core: 16 billion inst/sec

$\hookrightarrow 4-8$ per chip $\rightarrow 128$ billion inst/sec

\hookrightarrow floating point

2 or 3 floating point #s

call this gigaflops

balance 1 flop

floating point inst. per sec.

$\frac{1}{4}$ bytes

$\hookrightarrow 400$ GB/s

Memory technologies

Key tradeoffs in memory technologies:

Size (density)

compatibility \rightarrow how to integrate

speed/latency

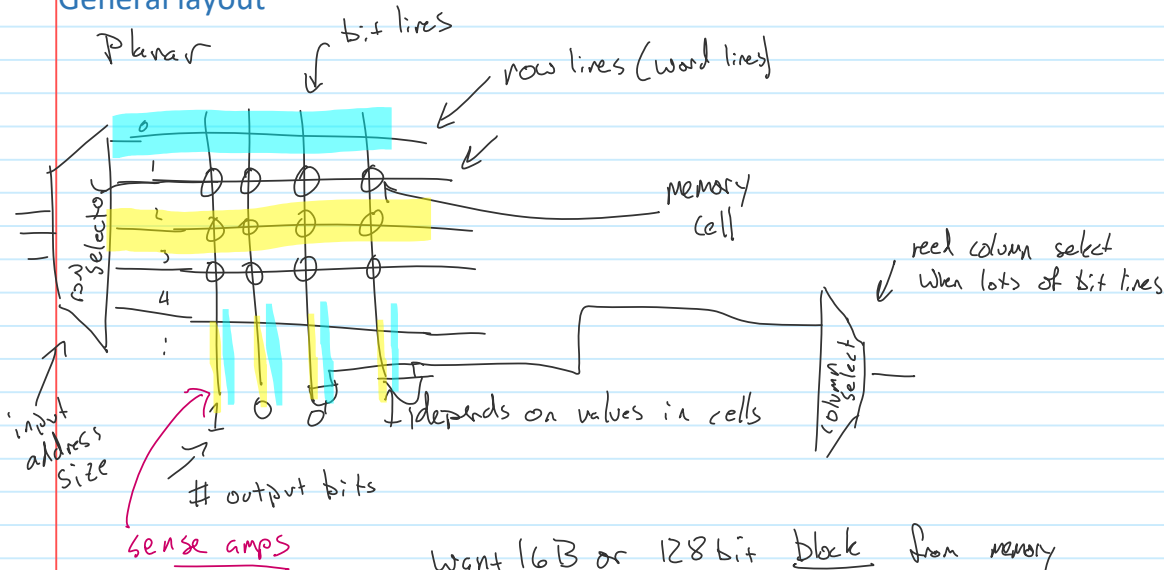
reliability (ECC)

cost

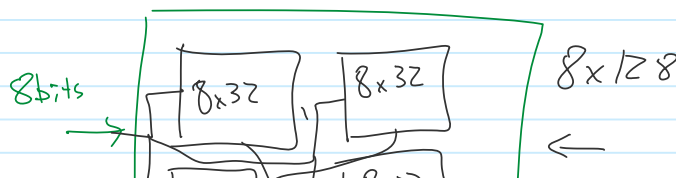
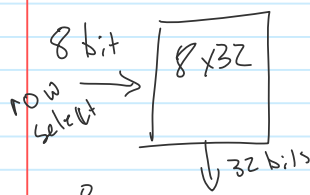
volatility (keeps data w/o power)

energy \rightarrow power
read/write \hookrightarrow static
dynamic

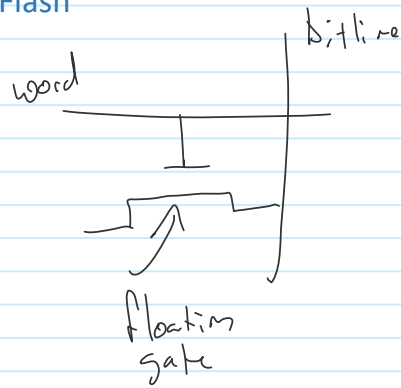
General layout



want 16B or 128 bit block from memory



Flash



Smaller than DRAM \rightarrow (no capacitor)
 \rightarrow store many bits per cell
 (MLC or TLC)

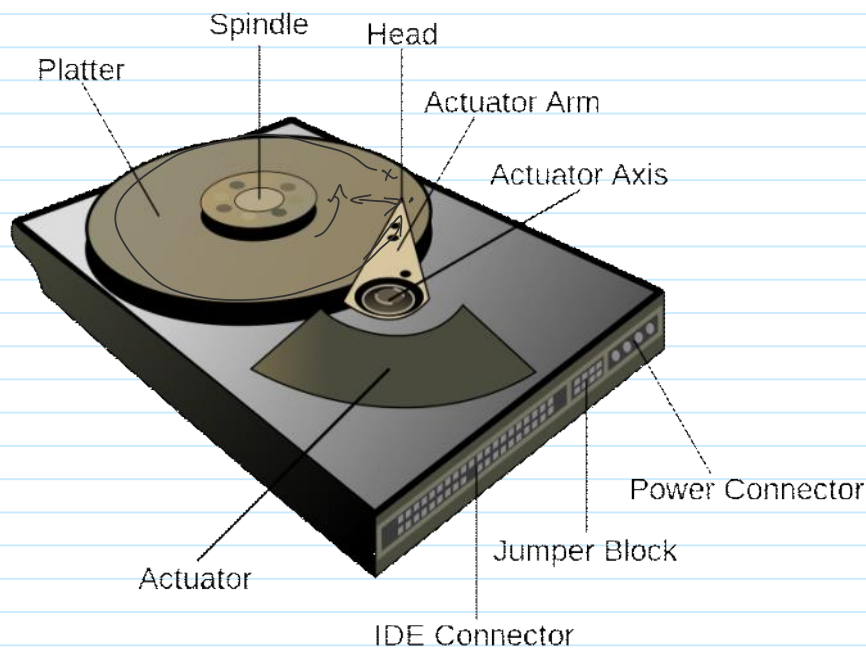
\rightarrow cheaper than DRAM

Non volatile

latency? $\sim 10-100\text{ ns}$ but must operate on large pages
 8-64KB

$\rightarrow 10\mu\text{s} +$

Magnetic disk



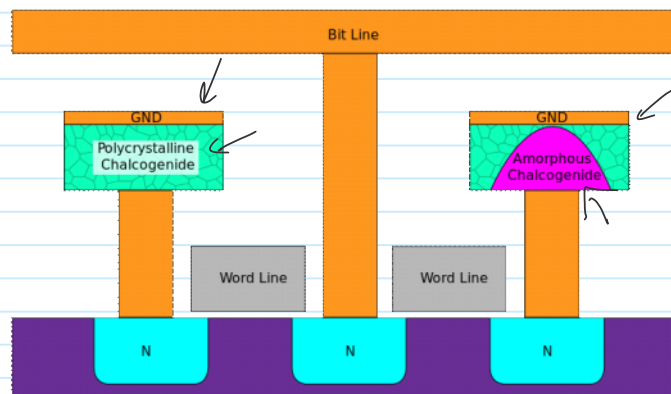
Size \rightarrow much smaller than flash

non volatile

latency \rightarrow much higher
 $10 + \text{ms}$

power \rightarrow high

Emerging technologies



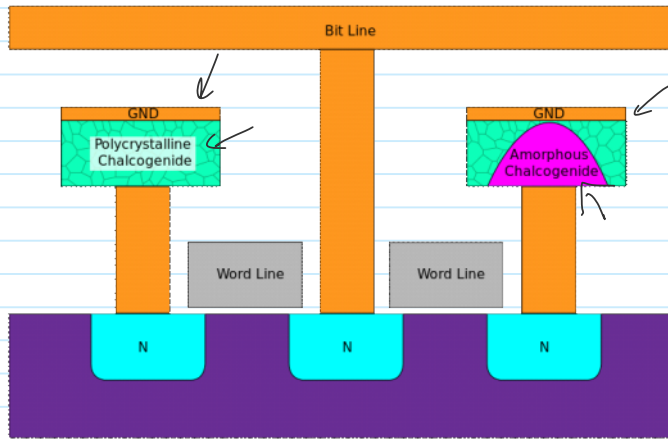
PCM phase change memory
 \rightarrow Intel's 3D-xpoint memory

Non volatile

theoretically very dense

latency \rightarrow slow to write $\sim 5-10\mu\text{s}$
 reads are fast $\rightarrow 1-10\text{ ns}$

power \rightarrow reading is low
 writing is high



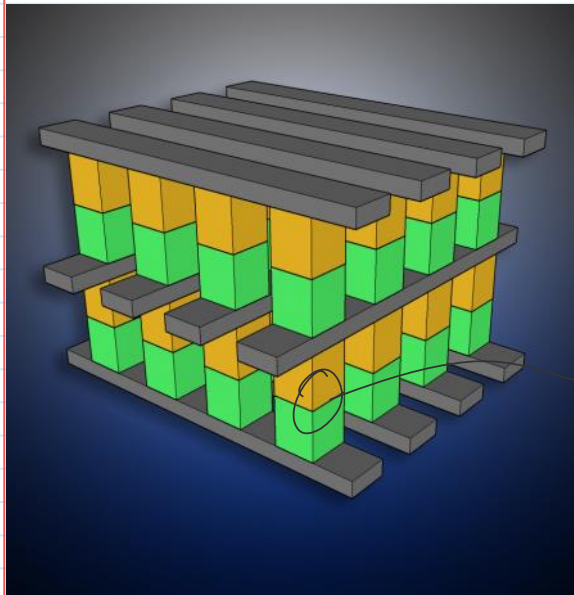
Intel's 3D Xpoint Memory

Non Volatile

theoretically very dense

latency \rightarrow Slow to write $\sim 5-10 \mu s$
reads are fast $\rightarrow 1-10 ns$

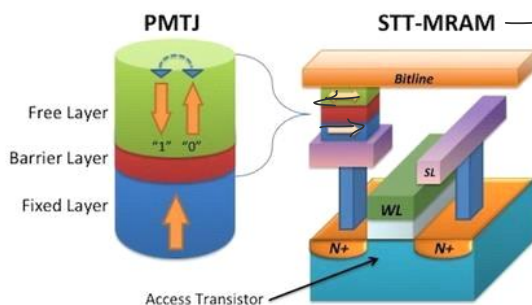
power \rightarrow reading is low
writing is high



Intel's 3D Xpoint

can be really dense by stacking in 3D

PCM



non volatile

\rightarrow Spin-transfer torque resistive RAM

low power and low latency