# Discussion 7: Hierarchy
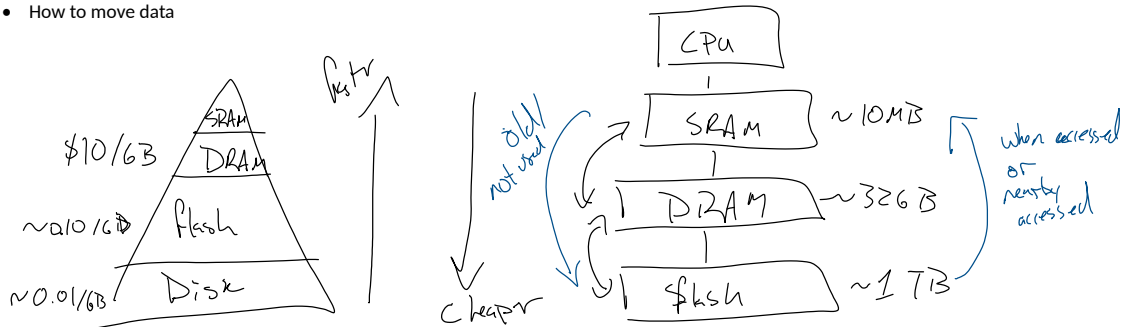
Thursday, February 21, 2019   2:34 PM

## Outline
- Hierarchy
- Locality
  - Temporal
  - Spatial
- How to move data

fastr ↑

$10/GB   SRAM
          DRAM
~$10/GB  flash
~0.01/GB  Disk

Cheapr

old/ not used

CPU
SRAM   ~10MB
DRAM   ~32GB
flash  ~1TB

when accessed or nearby accessed

how to decide when/what to move?
→ how often accessed
   **frequent** things close

   things close to eachother (sequentially)
   ↳ pre load / move together

temporal

locality

spatial

↳ "prediction"

CPU
SRAM    0-16kB
DRAM
        0-4GB

1) try to find data in SRAM
2) go to DRAM
3) put in SRAM

↳ hash table
   maps DRAM addr
   to SRAM addr

hardware cache

knights landing
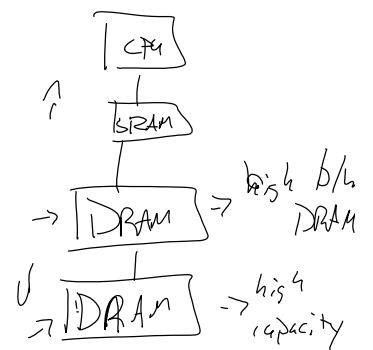
two pools of memory
Let programmer move data from
SRAM/DRAM and back

Scratch pad memory   or  → GPU
Software cache

More complex?
  hardware cache
  ↳ lots of metadata
    to track addresses present
    + where addresses' data is
                              ↳ tag + valid + state

CPU
SRAM
→ DRAM → high b/w DRAM
→ DRAM → high capacity

how to move data?

CPU  ──  SRAM  ── ── ── ── ──  DRAM  ── ── ──

data        addr

bus/wires ~64 bit + 32 bit
Parallel up to 512 bit
Short          → very fast → @ clock speed
On chip          Wide
                 low power → small wires

Off chip          limit
  ↳ through pins ↗
        onto fiberglass board
thick + long wires
  ↳ higher power
Clock is slower

narrower bus → pretty wide
        64 bits →

  ↳ convert electric
     signals to light
  ↳ expensive

On motherboard
  ↳ thick
Serial bus
  ↳ more lat
slow Package pin
low Bandwidth

Radio/ microwave interconnect
  ↳ low latency @ speed of light
  ↳ interference      ↳ convert to
  ↳ expensive  ↳ power    e/m waves

① Optical interconnect
  ↳ thicker waveguide
  ↳ low latency
  ↳ high bandwidth per "pin"
  ↳ incompatible w/ logic

3D die stacking
        through silicon
        vias (TSVs)
        copper wires like on chip
             wires
CPU
DRAM

wide busses    low power
high bandwidth     → heat        ↳ very expensive
low latency
               → design package!
               → new VLSI tools