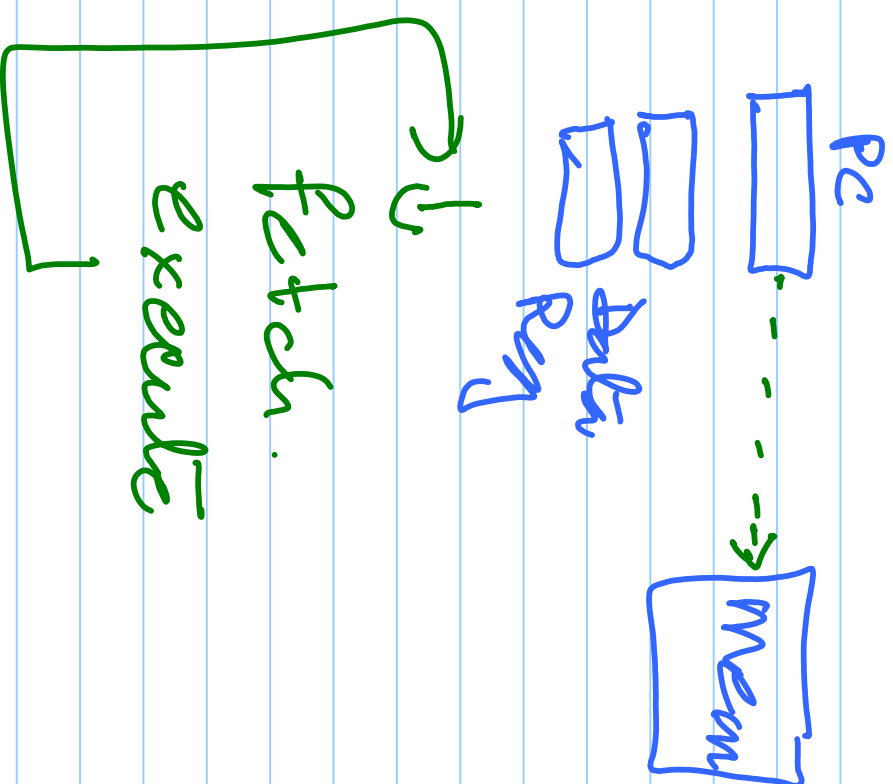


Single CPU

- single cycle CPU

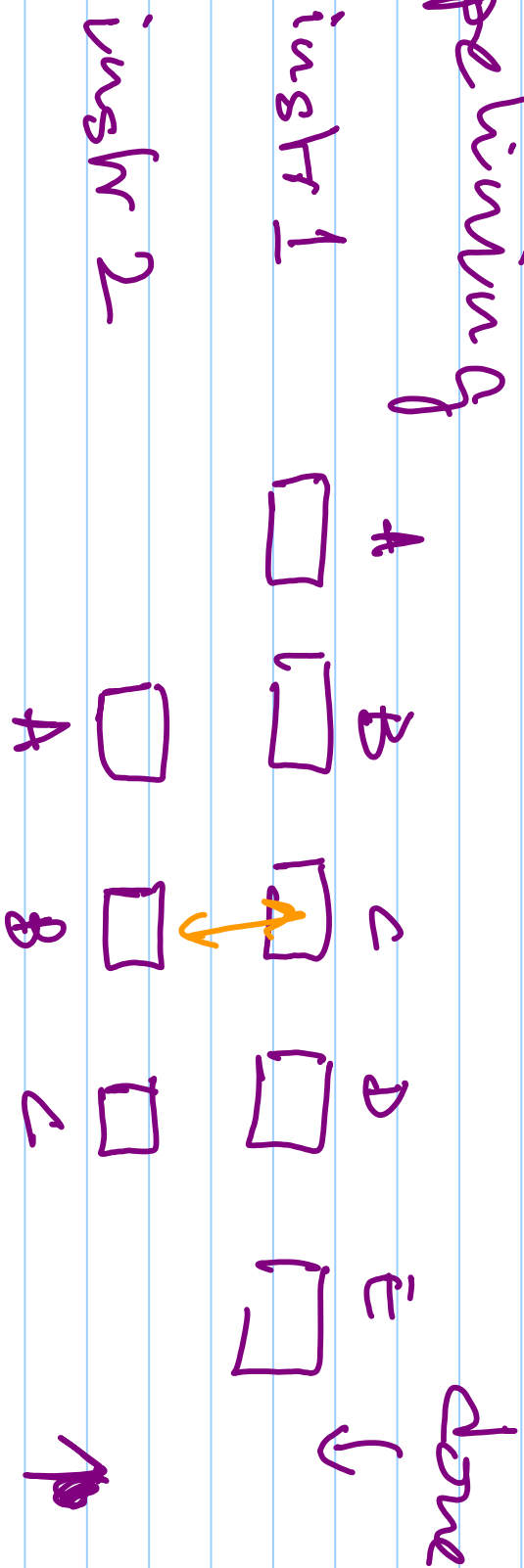
all software  
assume CPUs  
are single cycle



- Single cycle

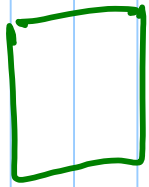
- cache → intermediate...

- pipelining



- hyperthr reading

Hyper

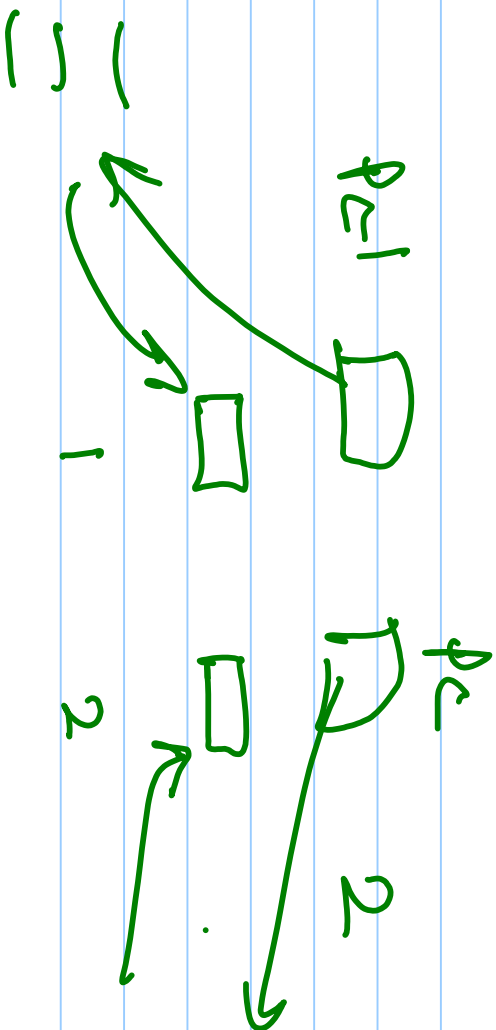


looks

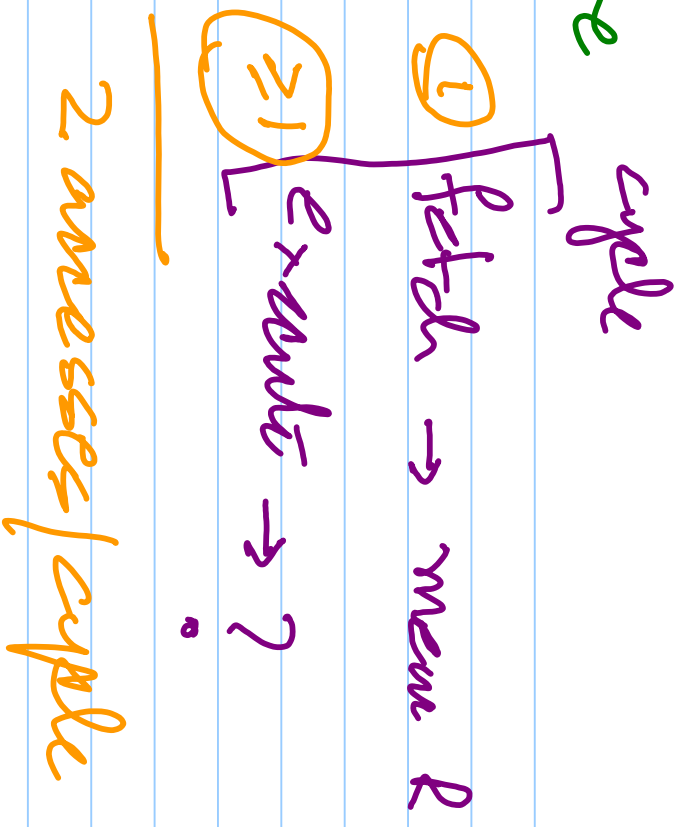
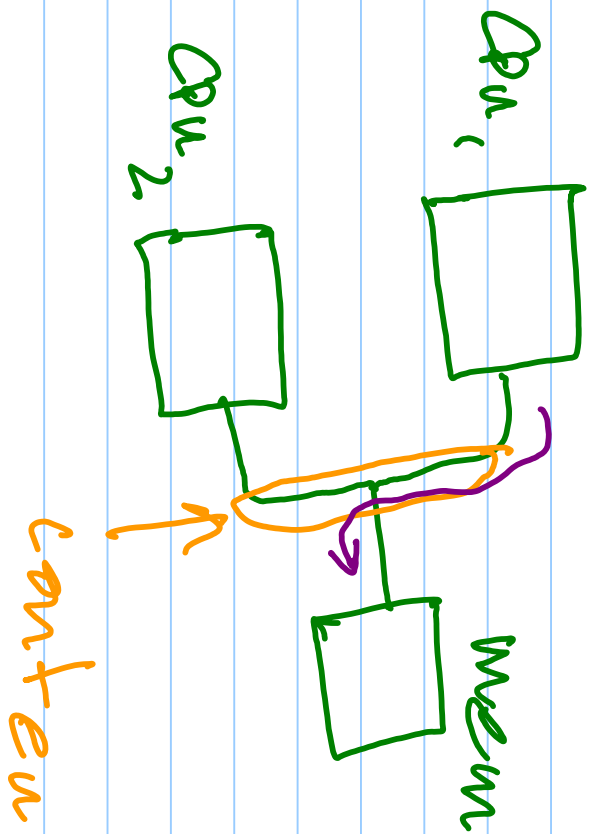
1 CPU

like

2 CPUs

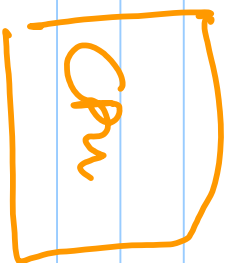


# VMA architecture share cycle

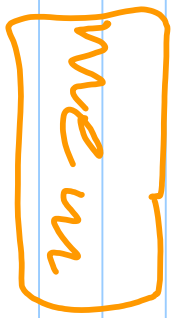


fastest

DRAM

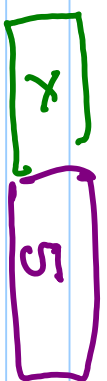


SRAM



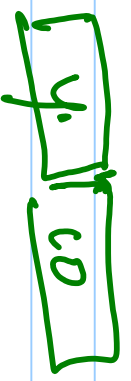
writes  $\leftrightarrow$   $x$

CPU

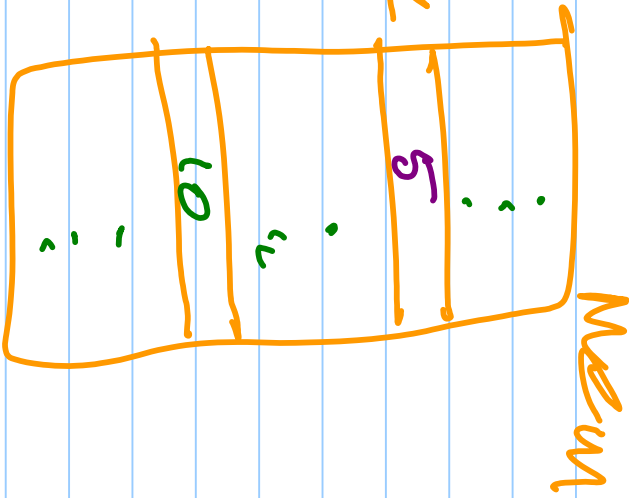


$x$

CPU reads  $y$



$y$



Memory

# Cache lookup

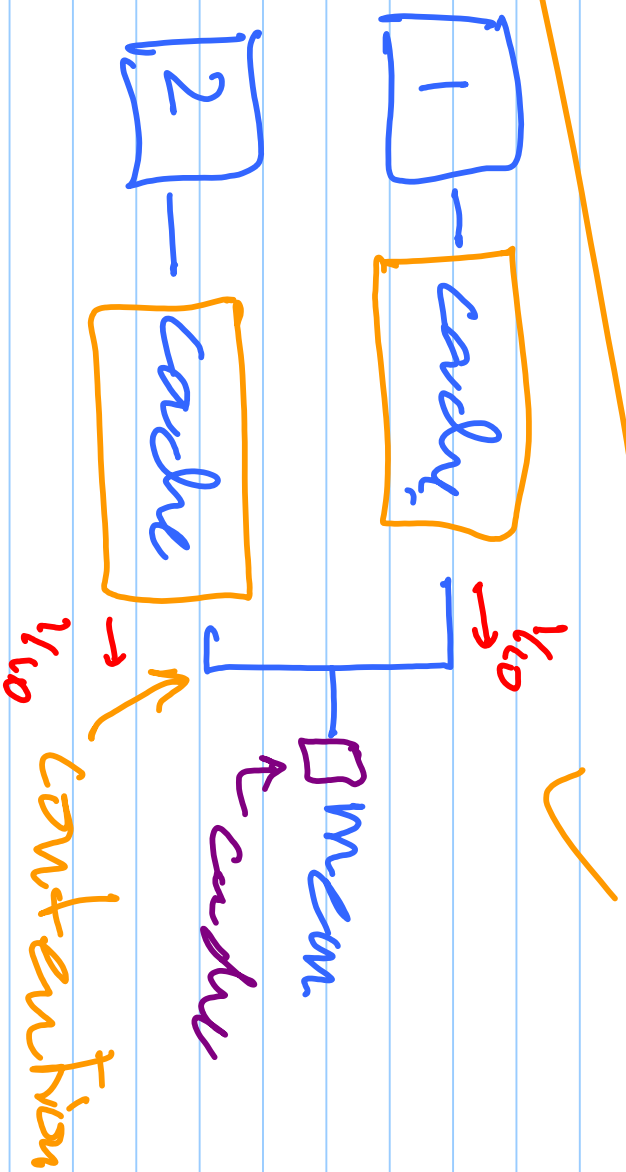
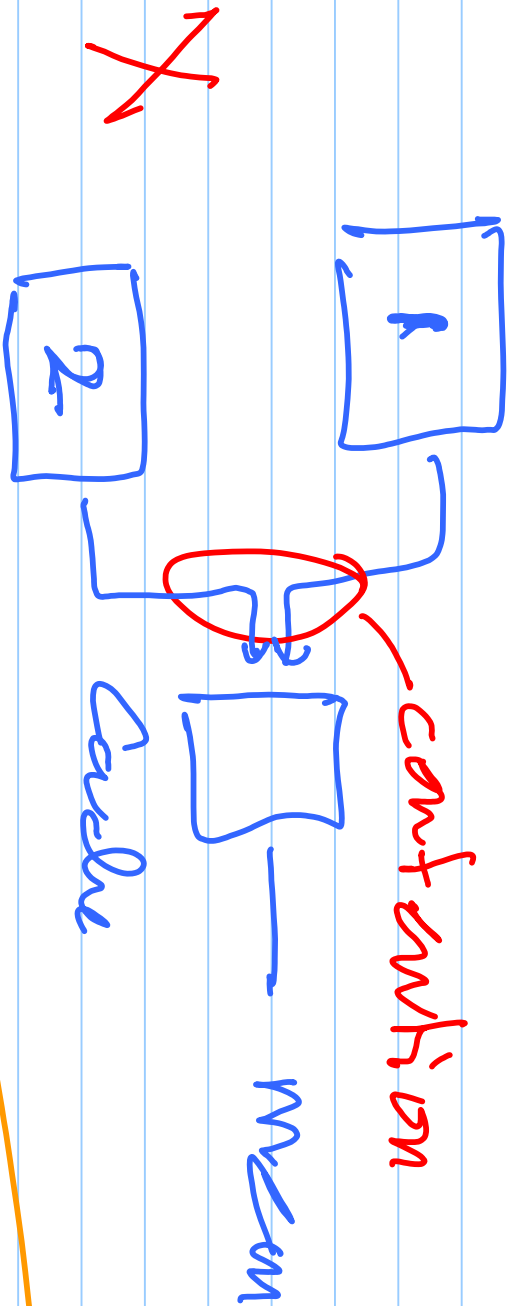
CPU → addr to cache

↓  
Search tags for  
addr

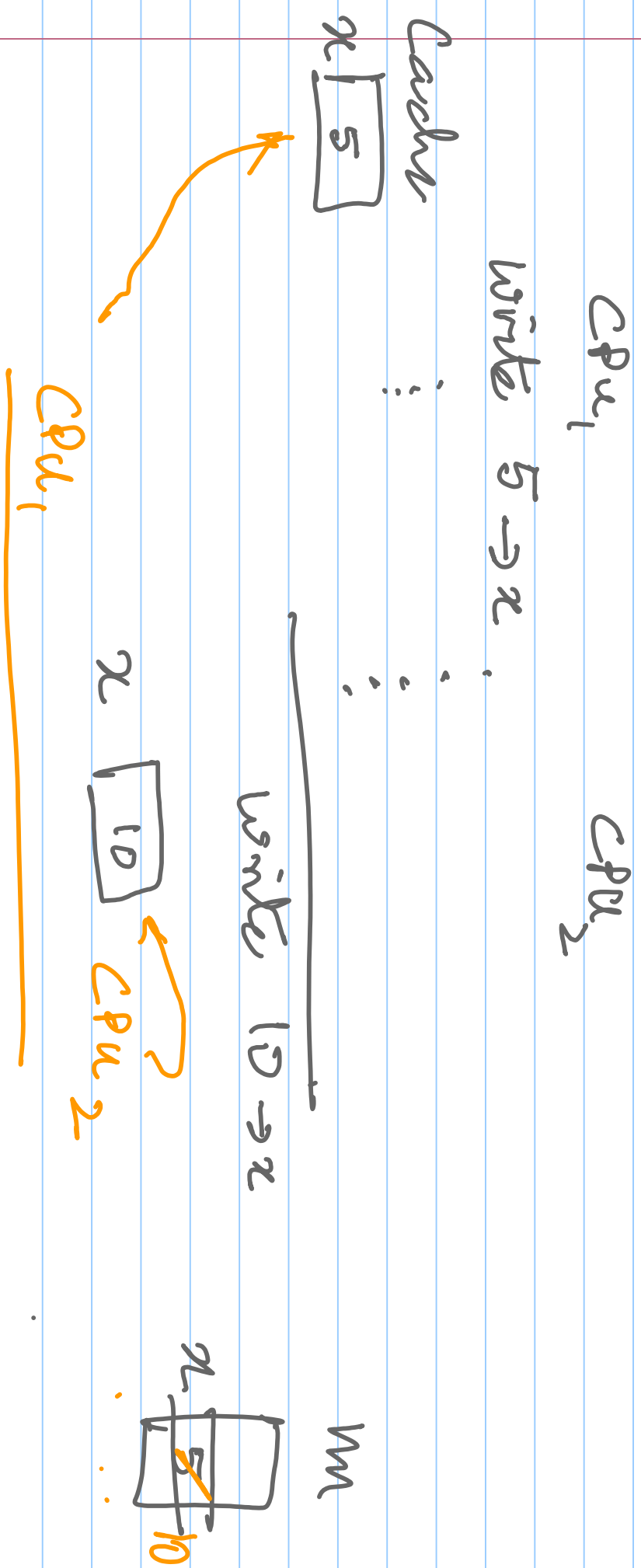
↳ exists → HIT

→ no " — MISS

# Multi CPU + cache



Problem  $\rightarrow$  Cache Coherency





non coherent cache

→ use programming to solve the problem

every read returns  
the most recent  
write

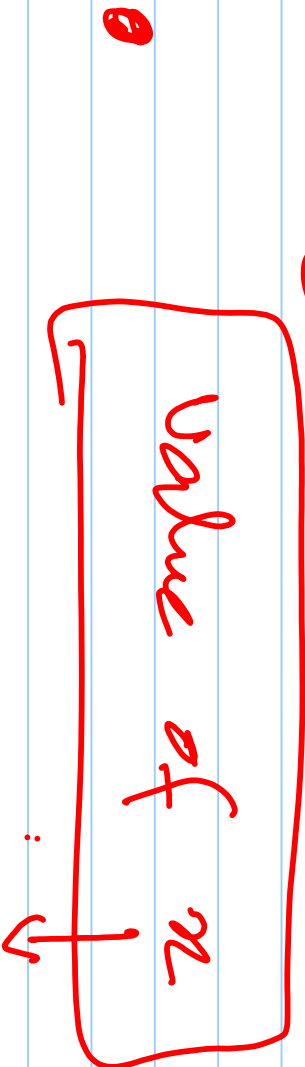
CPU<sub>1</sub>

CPU<sub>2</sub>

Write (S) to x  $\longleftrightarrow$  write (10) to x

---

read (x)  $\rightarrow$  10  $\leftarrow$  read (x)



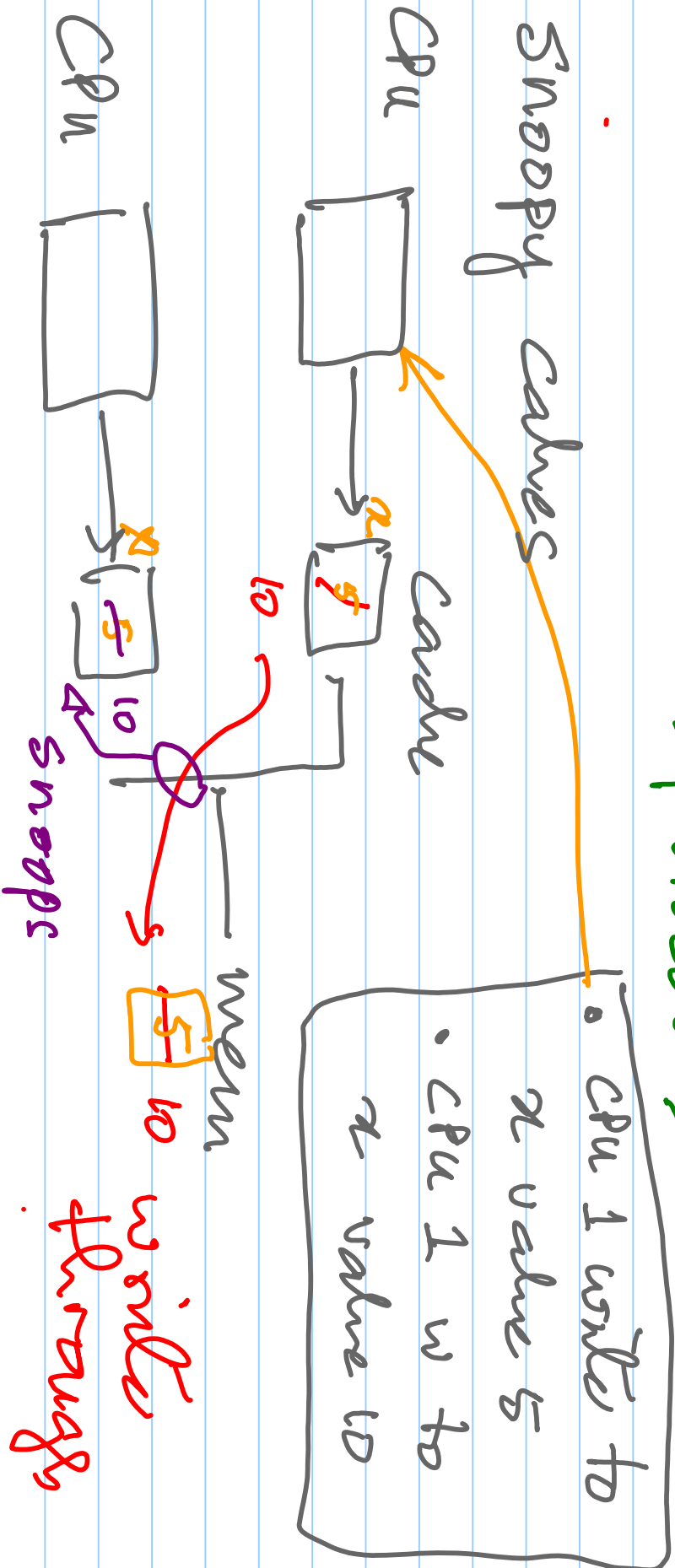
5 or 10

# Cohesive caches

(cache coherence)

→ protocol →

Snoopy caches



Coherent caches

↳ snoop

↓

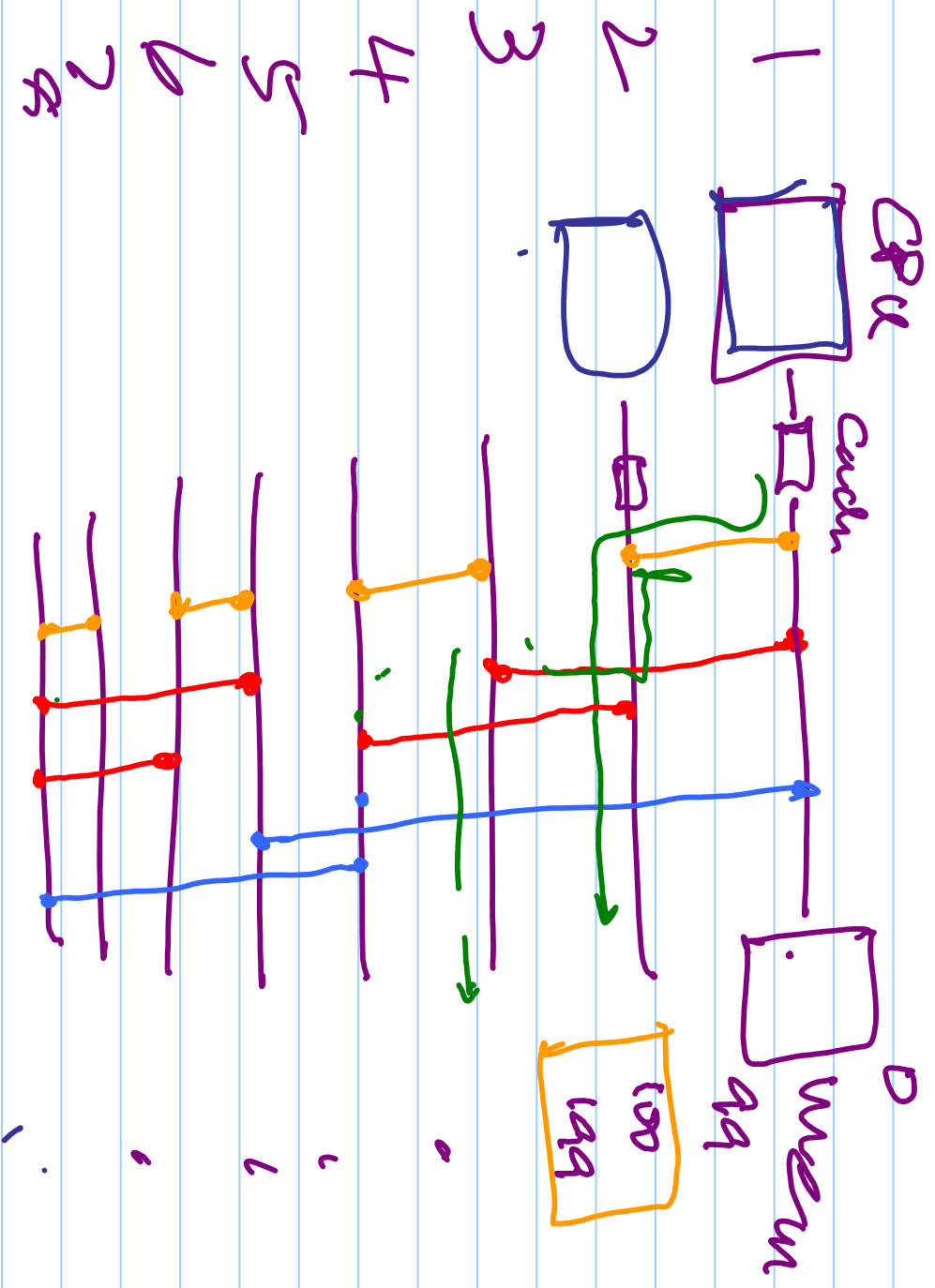
performance

degrades (but better than no cache)

↳ does not scale

# NUMA

NUMA X



hypercube

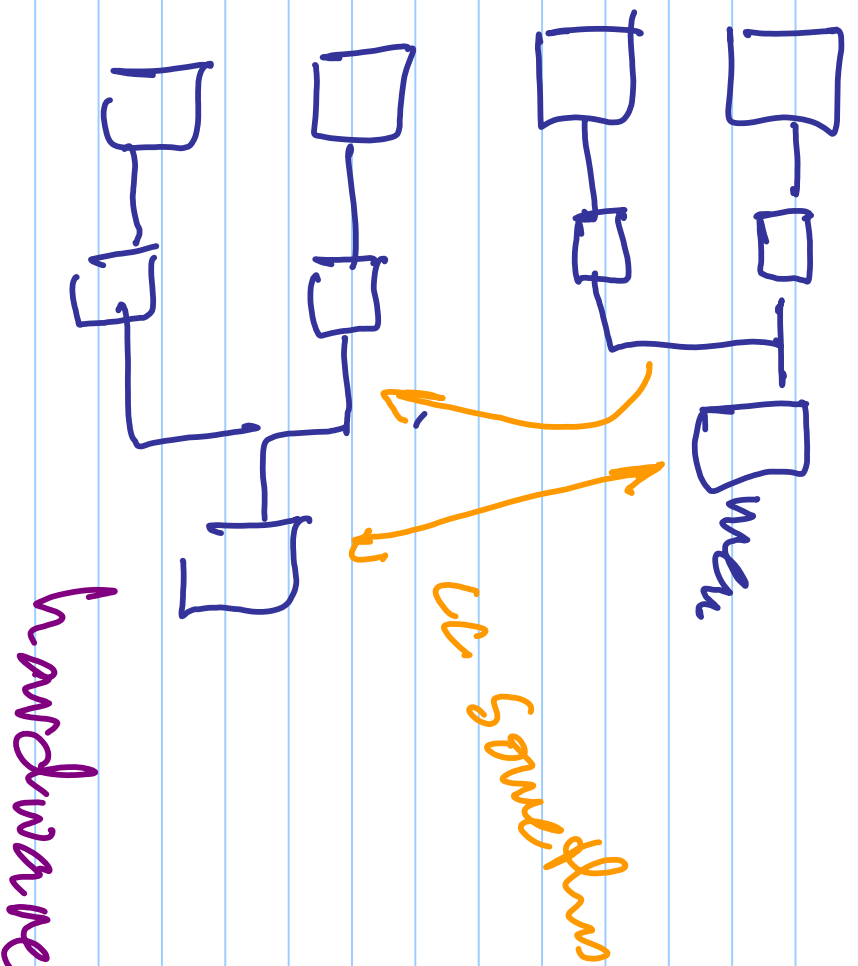
# NUMA

- Scalable
- high performance
- special appl
  - ↳ 1 appl / machine / time

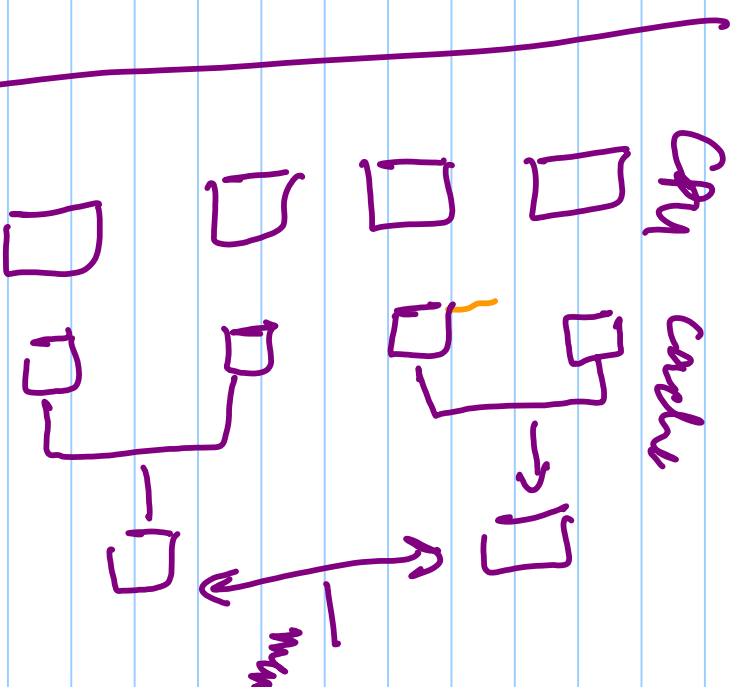
# CC-NUMA

Cache coherent NUMA

↳ NUMA that look like UMA



CC  
Sandy  
bridge



CPU cache

memory

linear bus