

Performance

- Can det from various metrics, but isn't easy
  - scientific sims. - fp performance
  - prog. dev. - int performance
  - commercial work - mem + i/o
- Response Time = minimize elapsed time for program (time-end-time-start)
  - ↳ important to computer system user
- Throughput = maximize completion rate = jobs per second.
  - ↳ important to computer manager.
- Without overlap in processing throughput =  $\frac{1}{\text{avg resp time}}$ , usually throughput >  $\frac{1}{\text{avg resp time}}$  otherwise.
- Performance for comp arch. is program run time (which we minimize)
  - performance =  $\frac{1}{\text{time}}$ , time = resp time, CPU run time, etc.
  - elapsed time = CPU exec time + I/O wait, comp arch focuses on exec time.
  - ↑ performance
    - ↳ faster cpus = ↓ resp time ↑ throughput
    - ↳ more cpus = ↑ throughput if resp time is possible.
  - Machine is  $n$  times faster than B if  $\text{perf}(A)/\text{perf}(B) = n$  and  $x\%$  faster if  $= 1 + \frac{x}{100}$ .
  - broken into instructions + cycles.
- Iron Law =  $\text{time}/\text{program} \approx \text{instrs/prog} \cdot \text{cycles/instr} \cdot \text{sec/cycle}$ 
  - sec/cycle (cycle time / clock time) - det by tech + L1 orgn.
  - cycles/instr (CPI) - det by ISA + CPU org., effective CPI not actual, overlap & CPI, average over instr.
  - instrs/program (instruction count) - instrs executed (not static), det. by program, compiler, + ISA
- Minimize overall time not just individual terms.
- Other metrics include MIPS + MFLOPs, but aren't as accurate,
  - (MFPS - mil float ops/sec), (MFLOPS - mil float ops per sec)
  - can be optimized but still worsen performance (is only CPI metric)
  - used in peak perf, MIPS is ok when instrs/prog is constant w/ same compiler + ISA.
- Benchmark programs are used to measure + predict workload performance
  - best case is running same set of progs everyday, but not realistic.
  - SPEC 2017 most recent - (System Performance Evaluation Cooperative) has 12 int + 17 float progs.
    - ↳ normalize runtime w/ a Gold Standard processor (SPEC ratio)
    - ↳ geo mean of the normalized times

Performance Contd

- How to average runtimes across programs when benchmarking?

- Arithmetic mean -  $\frac{\sum \text{time}_i}{n}$  ( $n$  programs)

x valid if programs run equally, otherwise use weight  
 $\frac{\sum \text{weight} \cdot \text{time}}{\sum \text{weight}}$

- Harmonic mean -  $\frac{1}{\frac{1}{\sum \text{rate}_i} / n}$  - for using rates & ending with rates.

- Geometric mean -  $\sqrt[n]{\prod \text{ratio}}$  - for using ratios

x not prop. to total time + usually mis predicts for  $\geq 3$  machines.

- Unnormalized raw runtimes are always best.

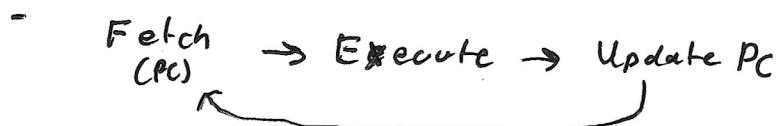
- Amdahl's Law - make the common case fast

$$\text{Speedup} = \frac{\text{time old}}{\text{time new}} = \frac{\text{unaff old} \cdot \text{aff old}}{\text{unaff new} + \text{aff new}} = \frac{1}{(1-f) + \frac{f}{s}} \quad (s \text{ factor of speedup, } f \text{ affected frac})$$

$$\lim_{s \rightarrow \infty} \frac{1}{1-f+\frac{f}{s}} = \frac{1}{1-f} \quad (\text{frac matters for max speedup})$$

Instructions

- words of a computer
- Instruction Set Architecture (ISA) is the vocabulary + defines the interface of computers.
- General purpose computer



- PC stores address of instruction + is updated to point to the next instr. to execute.

- Terms

- Opcode - specifies type of op
- operands - in + out data (must be registers)
- multiple instrs per C statement often

- Registers + ALU ops

- RISC-V uses 32 registers  $\#0 - \#31$  (32 bits/word)
- registers are useful for small storage for ops + fitting specifiers in instructions.
- memory load/store used to access + modify more data.
  - $\times$  address is byte address so word addressing ends addr. in two 0s in binary. (4 bytes per loc.)
  - $\times$  using load + store w/ offset is easy way for arrays
    - $lw \#5 0(\text{base addr})$  loads  $A[0]$
    - $sw \#5 4(\text{base addr})$  stores to  $A[0]$
- ALU ops ex: add, sub, mul, div, or, and, sll, etc
  - $\times$  can be done w/ immediates for  $\uparrow$  shift left (constants)

- Endian determines MSB @  $0x\dots 00$  or  $0x\dots 11$ . X86 is little endian. (Little end first (big end))

- Branches + Jumps - used to move around exec path

- branches compare two registers and either jump to an addr or continue regular execution.
- jumps move the addr to some other PC and sometimes stores the PC for  $\uparrow$  the next instr. in  $xt$  (jal or jsl),

- Procedure Calls
  - Initial responsibilities
    - × save reg\$, setup params., call procedure, get result, restore reg\$.
  - Process responsibilities
    - × do <sup>more reg\$</sup> func, set up results, restore more reg\$, return
  - jal is only hardware instr., rest is software.
- Stack - useful for calls
  - grows from large to small addr, w/ x2 being stack pointer
  - push x5 (add: x2, x2, -4, sw x5, 4(x2))
  - pop x5 (lw x5, 4(x2), add: x2, x2, 4)
- (at instruction level)
  - Software layers - possible bc pros. can be treated like data.
    - executable: datain → dataout
    - linker: executable file → executable in mem.
    - assembler: asm file → objfile
    - compiler: HLL file → asm file
    - editor: Commands → HLL file.
  - RISC-V machine language uses 32b instrs., while assembly instrs. are wordier.
  - Instruction format - reg\$ are 5b, opcode is 7b. (rs is source reg, rd is dest. reg.)
  - R-format - ~~func?~~, rs2, rs1, func3, rd, opcode
    - × used for three register ops. mostly simple operations
  - I-format - immediate [11:0], rs, funct3, rd, opcode
    - × used for immediate (constant) ops and loads. imm. is signed 11 bits
    - complement. imm. used for operand or load offset.
  - S-format - imm [11:5], rs2, rs1, funct3, imm [4:0], opcode
    - × used for stores
  - B-format - imm [12:10:5], rs2, rs1, funct3, imm [4:11:11], opcode
    - × used for branches, 0th bit not used bc 32b instrs. + 16b instrs. exist so only one bit not specified.
  - U-format - imm [31:12], rd, opcode
    - × used for lui, load upper imm.
  - UT-format - imm [20:10:11:11:11:12], rd, opcode
    - × used for jumps
  - instrs. structured weird so bits line up in every instr. Bits allow us to decide control signals to exec instr.

RISC-V

Other ISAs

- Addressing Modes

- register addressing
- base addressing - base addr + offset
- Immediate addressing - immediate values
- PC relative addressing - address based on PC + immediate
- Pseudo direct addressing - address in instruction.
- Indexed addressing - add two ress base + index
- Indirect addressing -  $M[M[\text{addr}]]$
- Auto incr/decr - add data size
- Auto update - like displacement/offset but update reg.

- Choose ISA based on hardware implementation + optimized compilation.

- $\downarrow$  mem uses dense fast ISAs usually, w/ complex + irregular instrs.
- $\uparrow$  speed uses pipelines (v important) so simple + regular is important.

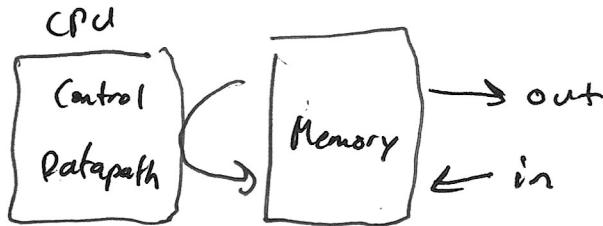
- Intel 80386 (x86) ISA is a decoding nightmare w/ 1-17 byte instrs, Prefix + postfixes, + weird format, but good core instrs.

- all instrs have to work in x86 tho.
- Intel chips translate into RISCy ops for execution units due to backwards compatibility + complex decodings
  - X helps to speed up execution + parallelism.

- More powerful/complex instr. doesn't mean faster.

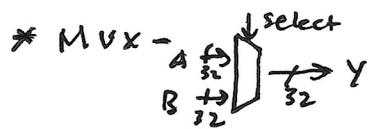
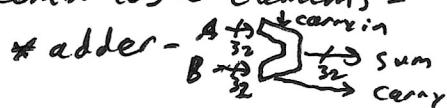
Single Cycle

- CPU/processor -

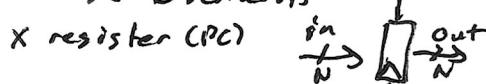


- Data path for instructions (single cycle) ~ composed w/  
mux, FF, gates, ALU, + reg file.

- Comb. logic elements -

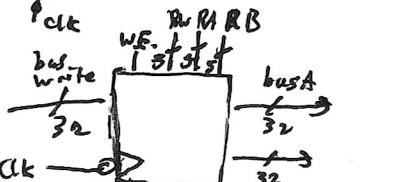


- Storage Elements



- \* Register file

- 32ress,  
2 read ports  
1 write port.

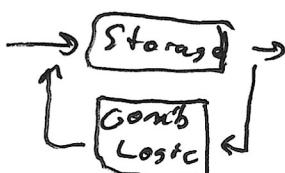


- \* Memory

- 1 in bus, 1 out bus,  
not bidirectional



- Computer:



- Single Cycle data path needs to do ~~one instr~~ in 1 long clk cycle.

- fetch, decode/read operands, execute, memory, writeback

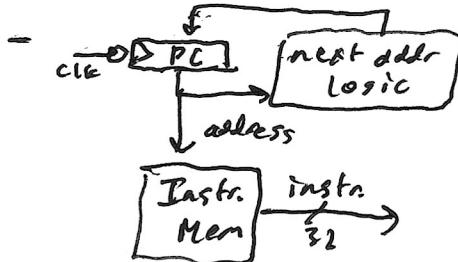
- RTL is useful to represent steps + dp. elements.

Single Cycle Contd

- Register Transfer Language (RTL) is nice short hand similar to pseudocode which helps describe what instructions do.

- ex.  $R[rd] \leftarrow R[rs1] + R[rs2]$        $PC \leftarrow PC + 4$  (add)

- Fetching instructions - fetch & update PC at end of cycle



- next addr logic w/o jumps or branches:



- ALU instructions - utilizes the reg file & ALU (Reg write + ALU op)

- Logical ops. w/ immediates - utilizes reg file & ALU & (Reg write + ALU op)

- mux to switch alu input to immediate unit & mux + ALU sel

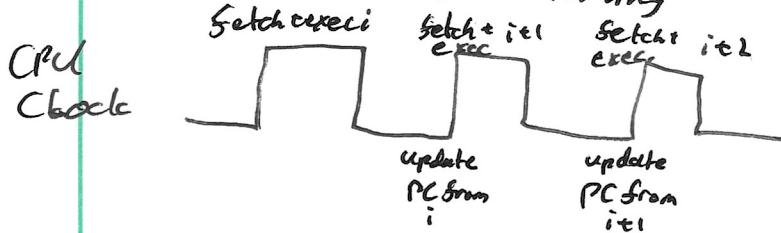
- Immediate unit to extract & sign extend value.

- Load instructions - uses reg file & ALU to calculate base addr. + offset, data memory, & a mux to choose between memory & alu output for writing data to reg file

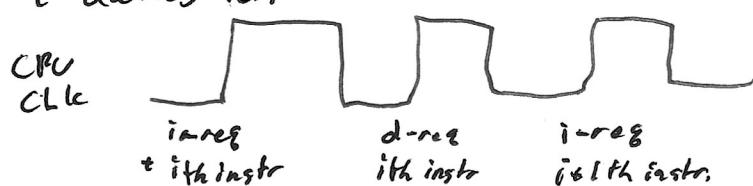
- Store instructions - uses reg file & ALU to calc. base addr (ALU op & ALU sel + offset, data mem, and a direct data in line from RS2.)

~~Ch~~

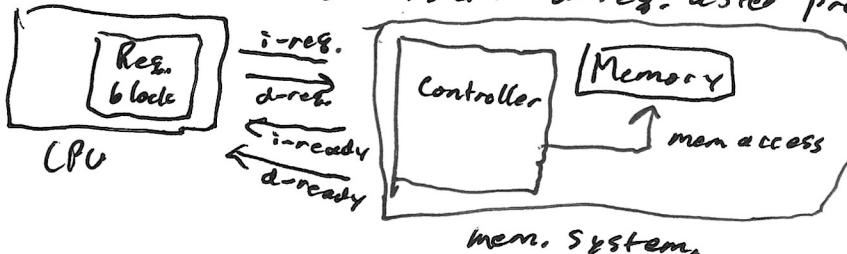
- Ideal Single cycle timing



- Some FPGAs require multiple clocks for mem. access for arbitration between Instr. fetch & data fetch



Single Cycle Control

- Memory Arbitration
  - Decides on what mem req. to process, as i-reg is generated in all instrs. including doregs.
  - Prioritizes d-reqs > i-reqs (i-reg will be fetching next instr. while d-req. is for current instr.)  
X Will have to deassert d-reg. after processed
  - 
  - CPU mem interface is seq. while controller is combinational.
  - Need to latch instr. while waiting for mem. to prevent instr. vanishing.
  - Needs to work at varying mem. latency (mem access time unknown in # of cycles). Allows to check for timing independence.
- For jump + branch instrs. additional muxes + adders + connections are made to update PC and sometimes registers.
- Control signals control datapath units for each instr based on opcodes, functs, and other parameters.
- Cycle time or max CPU clk speed depends on the design.
  - Time to compute next state values
    - X propagation delay  $Clk \rightarrow Q$ : (58)
    - X comb. logic delay
    - X setup requirements for ff.
  - Often times, load instr. have longest critical paths in single cycle due to data mem access.
    - X Leads to very long cycle time + slow processing.

Pipelining

- 1st of 4 big ideas of modern computer systems.
- Pipelining benefits (improving performance) - ↑ throughput, ↓ latency
  - Breaking instructions into stages + overlapping speeds up cycle time while keeping CPI  $\approx 1$ .
- Common Pipeline Stages
  - Ifetch - instruction fetch (from mem)
  - Reg/dec - Reg. read + instr. decode
  - Exec - compute/calc mem. addr
  - Mem - read data from memory
  - Wr - Write data to reg. file.
- Pipeline speedup
  - Ideal speedup = # of pipeline stages (equally balanced)
  - Max speedup  $\leq$  # of stages, Speedup  $\leq \frac{\text{Time for unpipelined op}}{\text{Time for longest stage}}$
  - Async / self-timed pipelines don't follow this but aren't common due to difficulty for high perf complex CPUs.
- Pipelining is implemented using latches (ffs) which keep the instrs. and transfer them between stages.
  - ↓ prog. latency by ↑ instr. throughput but keeping or increasing instr. latency (due to things like latch overhead)
  - Associate resources w/ states
    - Resources not necce indivisible/full-cycle (Write to res first half, read second half)
    - Carry necessary control signals along.
- Observations
  - Each hardware block only used once per instr. & must be used at same stage every time.
    - X R-types instr. don't use mem. so just do nothing during mem. stage.
  - Control Signals - generate once & pass along unconsumed signals through each stage

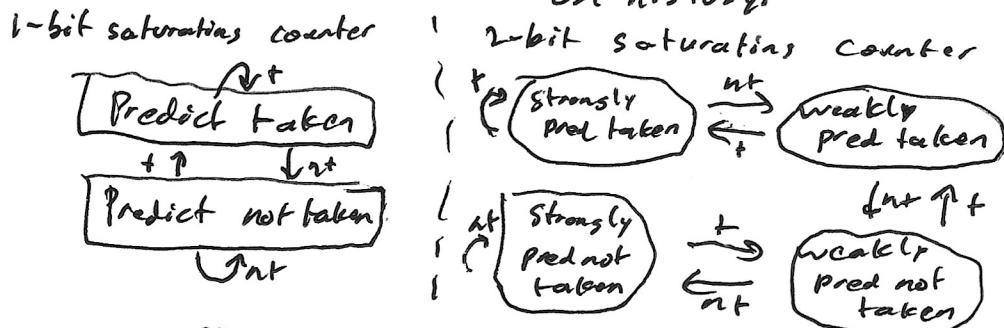
Pipelining contd

- Pipeline hazards - issues arising from the pipeline system.
  - Must be detected to address
  - Always can be solved by waiting, w/ stalls, bubbles, ~~stalls~~
    - ✗ ↑ CPI, ↓ speedup
- Structural Hazard - 2 instructions need same hardware at same time
  - Single memory, solve similar to S.t., process d-reg. first.
  - Write half cycles, read next half cycle.
- Data Hazards - data not ready
  - dependencies backwards in time, true dependence is forward in time
  - very common, so stalling slows down majorly
  - utilize a forwarding unit to forward data from future latches back to exec stage + muxes
  - RAW (read-after-write) (forwarding), WAR (allow in final stage), WAR (fetch operands early in pipe)
    - ✗ RAW only real hazard.
    - ✗ RAW detected from pending writes + operands having same reg.
    - ✗ WB stage no hazard bc write in 1st half cycle then read.
  - Loads still introduce RAW bc value not avail. till mem stages need to stall. earlier stages
    - ✗ True backward dependency
- Terms
  - Hazard - true backward dependence in time (can't be solved by forwarding)
  - Dependencies - in code, cannot be solved away (not always a hazard)
- Control Flow Hazards - not known which instr. to fetch.
  - branch instrs. + jump instr. (jump is calculated but ~~can't~~ may be avoided)
  - somewhat common so bubbles <sup>usually too much overhead</sup> lead to major delays.
  - can solve with delay slots w/ compiler, putting indep. instrs. in delay slots after branch
    - ✗ low use as long pipelines don't effectively use + indep. instrs. to be placed aren't very common.
    - jumps will almost always trigger this apart from advanced solutions.

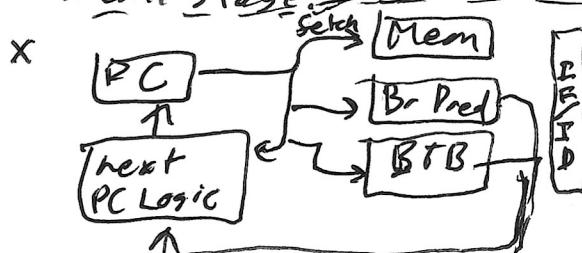
Pipelining ~~contd~~

- Branch Prediction - combats c.f. hazards by choosing a branch to follow.

- Easiest - predict not taken always (PC+4 easy to calc)
  - x Continues execution seamlessly.
- Need to flush and fetch correct instr. if incorrect (same delay as bubble so no issue in guessing)
- Dynamic branch prediction - adjusts prediction based on history.



- x works off recent history
- x store state of predictor in branch predict table w/ PC used for index + stored as tag (replacement occurs if collision).
- x store prior targets of branches of taken in branch target buffer in a similar manner to allow access to <sup>branch</sup> PC immediately
- x update BTB + BPT when instr. is resolved in mem. stage



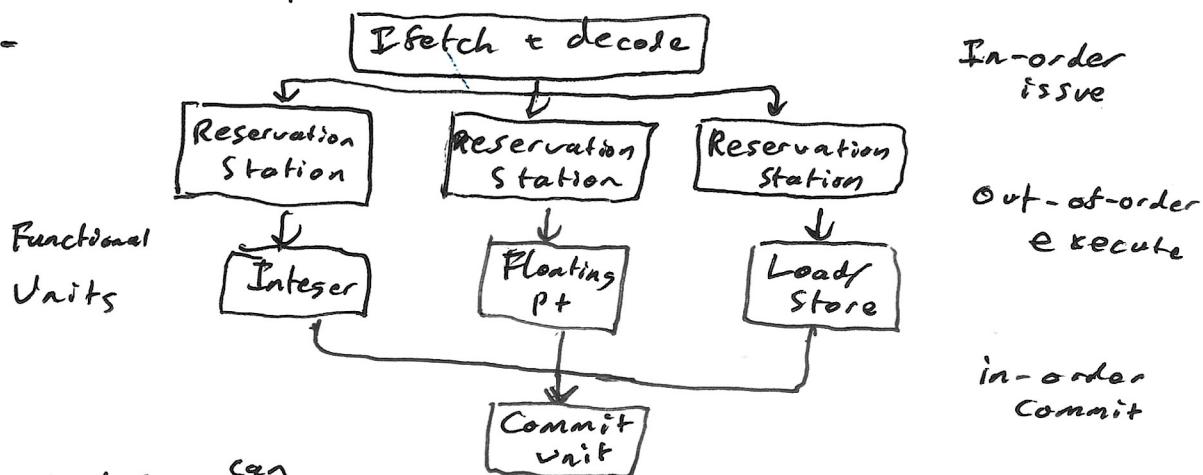
- When Calc'ng CPI do  $CPI = \text{ideal} + \text{frac. stall t...}$ 
  - Remember to calc branch & branch pred error.

Pipelineing Contd

- Exception - unprogrammed surprise function call (external)
  - handled by OS, must record the address of the offending instruction
  - returns control to user by saving + restoring user state
  - lets OS yet in + out without breaking program.
- Interrupts - caused by external events (e.g. I/O)
  - async to prog. exec.
  - can be handled between instrs.
  - suspend user prog., handle, and resume user prog.
- Traps - caused by internal events (~~user errors~~)
  - caused by exceptional conditions, errors, or faults.
  - sync. to prog execution
  - cond. must be remedied by ~~other~~ handler of prog.
  - instr. may be retried/simulated + prog. continued, or prog. may be aborted.
- OS wouldn't work w/o exceptions + interrupts
  - uncommon case but must implement
- RISC-V defines instructions having no effect if it causes an exception.
- Pipeline Considerations for exceptions
  - no instr. after the excepting instr. may execute
  - all instr. before excepting instr. must execute
  - excep. handler software saves + restores reg. and other state.
- RISC-V Exceptions - all jump to same handler code
  - we consider illegal instrs., arith. overflow, + hardware malfunction.
  - Hardware behavior
    - X save PC of offending instr.
    - X use special reg. SREG
    - X set SCAUSE reg appropriately: ELL=0, OVF=1
    - X jump to handler at fixed addr.

Pipelining cont'd

- Datapath modifications for exceptions
  - Hardware malfunction - detected in EX, squash (convert to NOP) + earlier stages
  - Illegal instruction - detected in ID, squash (convert to NOP) + ~~in~~ <sup>earlier</sup> stages
  - Need exception sorting when overflow followed by illegal exception (causes simultaneous exceptions)
- Superscalar Processor - multiple scalar ops/instrs in one cycle.
  - duplicate hardware for alu, imm gen, more lines out of register etc. for multi scalar ops per cycle.
  - static multi-issue processor
- Dynamic Scheduling
  - helps avoid/reduce effect of hazard
  - execute indep instrs. while waiting on results of <sup>Prior</sup> instrs.



- instrs. can execute when operands are ready
- instrs. "commit" when all prior instrs. have committed

Caches

- 2nd of 4 big ideas of modern computer systems
- CPU is majorly slowed by mem. system
  - real mem can take 300+ cycles.
  - large mem is slow + cheap, small mem. is fast + expensive.
  - Solved by mem hierarchy of progressively smaller + faster mems.
    - × provides capacity of large mem. + speed of small mem.
- Locality - programs access small portions of address space before moving on.
  - mem. access, data structures, instructions, are all often benefited by taking advantage of their locality.
  - Temporal Locality - Locality in time, keep recently accessed data close to processor.
  - Spatial Locality - Locality in space, move blocks of data to upper levels instead of single entries.
  - Move data to cache on access for temp. locality and move a block instead of a single word for spatial locality into fast mem. / cache
    - × As locality changes move new data in + replace old data in cache.
- Terms
  - Hit - data is in upper level (Cache)
    - × Hit Rate - fraction of mem accesses to upper level
    - × Hit Time - time to access upper level: ~~hit~~
  - Miss - data needs to be retrieved from block in lower level.
    - × Miss Rate - 1 - Hit Rate
    - × Miss Penalty - time to fetch block from lower level + time to deliver block to CPU
    - Hit Time << Miss Penalty.
- Caches take advantage of Amdahl's Law, as cache-hits through locality are the common case.
  - Reduces average memory access time (AMAT)
    - ×  $AMAT = \text{hit time} + \text{miss rate} * \text{miss penalty}$

Caches <sup>coated</sup>

- Managing the memory hierarchy
  - Register  $\leftrightarrow$  Main memory - managed by compiler (sometimes programmer)
    - $\times$  Word granularity, but sw ties mem loc. to registers (alloc), temp. vars to spill mem. when needed.
  - Disk  $\leftrightarrow$  Memory - managed by OS for mem.
    - $\times$  Automatic & transparent to user
    - $\times$  Managed by programmer for file row
    - $\times$  Virtual memory & illusion of largemem. w/ protection.
  - Cache  $\leftrightarrow$  Memory - managed by hardware for speed
    - $\times$  auto managed
    - $\times$  Block granularity for spatial locality, retaining recently used blocks for temporal locality.
- Cache - Hardware Hash Table
  - Tag - indicates which address is in the slot
  - Data - block of data
  - Valid - whether the data in the entry is valid.
  - Block/Cache Frame - each slot in the cache.
- Hit/Miss detection
  - tag matching  $\Rightarrow$  hit & returning data to processor
  - else, miss, evict a ~~old~~ block, set block from mem into Cache.
- Address: Tag, index, block offset, byte offset
  - Byte offset - bits to index bytes in each line in the cache block
    - $\times$  ex. 4 byte line = 2 b byte offset [1:0]
  - Block offset - bits to index lines in each block
    - $\times$  ex. 2 lines per block = 1 b block offset [2]
  - Index - bits to index blocks in cache
    - $\times$  ex. 64 blocks in cache = 6 b index [8:3]
  - Tag - leftover bits to store in cache to rep. address.
  - Block offset can include byte offset, ~~but not~~ Since addr. are byte addr.
  - More detailed expl. later.

Cache  
Frame

Caches ~~contd~~

- Block Placement Schemes
  - Direct mapped - many  $\rightarrow$  one - block can only be in one frame
  - Fully associative - many  $\rightarrow$  any - block can be in any frame
  - N-way associative - many  $\rightarrow$  N - block can be in one of N-frame
    - $\times$  N tags in set searched in parallel.
    - $\times$  Index selects set, tags are compared, + data is selected.
- Terms.
  - Cache Block - aka cache line
  - way - one of the assoc. ways in a set
  - set - set of all the ways/blocks in a set
  - index - selects the set
  - frame - physical loc. for 1 cache block.
- Cache Organization
  - Frames = Cache size / block size
  - ~~# sets~~ = Frames / ways
  - Block-offset bits =  $\lg(\text{Block size})$
  - Index bits =  $\lg(\# \text{sets})$
  - Tag bits = addr. bits - index bits - block-offset bits
- Set-assoc. cache adds another mux + have to det hit/miss early.
- Associativity Spectrum
  - Conflict misses  $\downarrow$  with higher assoc. due to cache mapping
  - $P_{\text{assoc.}} = P_{\text{tag bits}}, f_{\text{idx}}^{\text{idx bits}}, T_{\text{comparator size}}$ 
    - $\times$  assoc. search is complex.
  - Direct Mapped  $\xrightarrow{\text{missrate} \downarrow}$  Fully Assoc.
 
$$\xleftarrow{\text{Complexity} \downarrow}$$
  - A little assoc. goes a long way in  $\downarrow$  miss rate.

Caches contd

## • Miss Types

- Compulsory - cold start / process migration - first access to every block  
x can't do much, insignificant if running tons of insts.
- Conflict - collision - multiple mem. loc. mapped to same cache loc.  
x Can T cache size / associativity
- Capacity - cache can't contain all blocks accessed by the program  
x can T cache size
- Coherence - invalidations - other process updates mem. (multicore / I/O)
- 

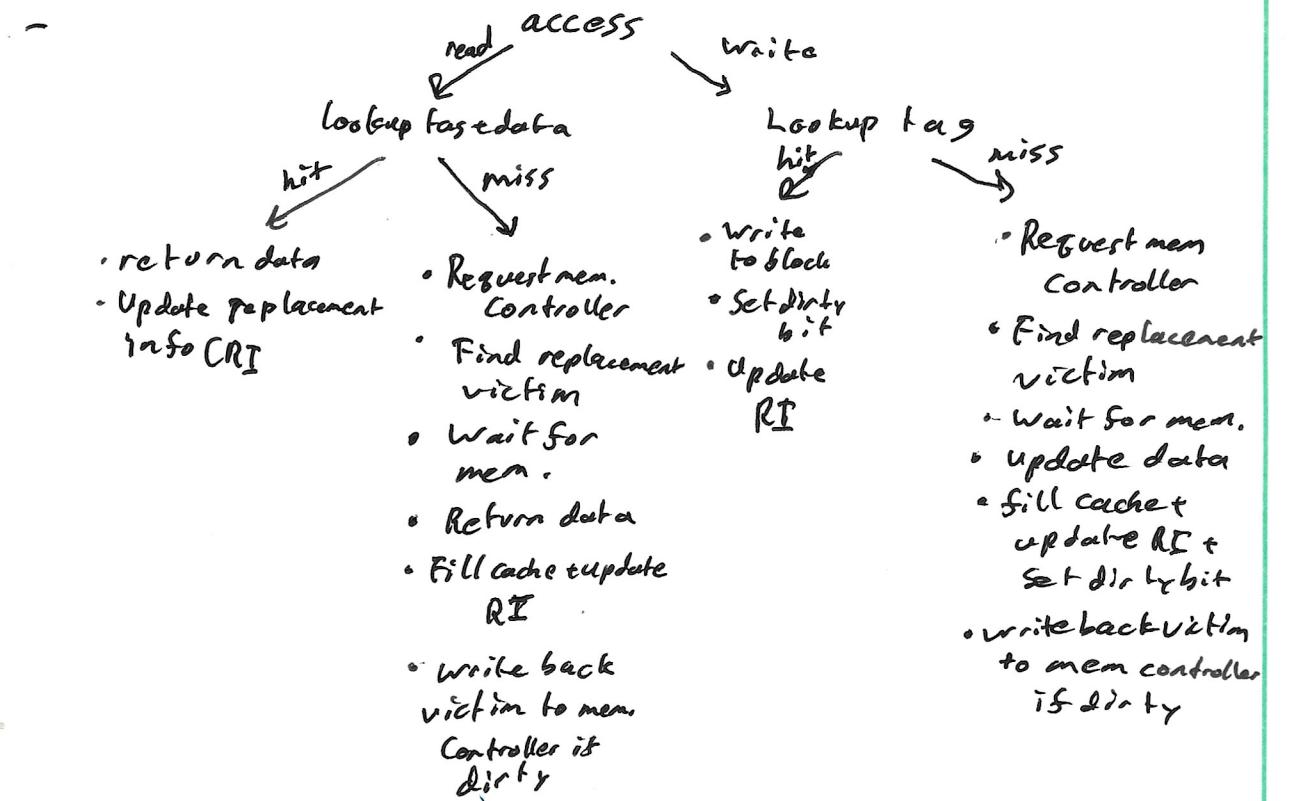
Cache Size	Direct Mapped	N-way set assoc.	Fully assoc.
Compulsory Miss	Big	Medium	Small
Conflict Miss	Same	Same	Same
Capacity Miss	High	Medium	0
Coherence Miss	Low	Medium	High
	Same	Same	Same

## • Block Replacement Policies

- Direct Mapped is easy (have to replace block fills up spot)
- Assoc. Caches
  - x Random - simple
  - x FIFO - not usually good
  - x LRU (least recently used) - smart
  - x NRU (not recently used) - not bad.
- Handling Writes
  - Could queue up another store while handling the current one, as stores don't write to register.
  - Write-through on writes - write to cache + memory (write buffer)
    - x simple but of mem traffic.
    - x write buffer can be fifo from the store queue
    - x using 12 cache ↓ delays to overflow
  - Write-back → write to cache and write to mem. when block is replaced.
    - x only one write to mem. (write hits)
    - x reads can cause writes (evicting blocks)
    - x multicore coherence (non-intuitive)

Caches Contd

## • Cache operation (WB)



- Separate Mem controller (controls RAM) & cache controller (controls cache operations + states)
- Mem. controller only acts upon C.C. requests.
- need to latch reqs from C.C. to m.c. to prevent large critical path from addr  $\rightarrow$  cache  $\rightarrow$  mem.  $\rightarrow$  cache  $\rightarrow$  pipeline.
  - $\times$  slows clock even further than single cycle often.
  - $\times$  latch all req + data between controllers
    - o data from mem + req. from cache primarily

## • Block Size Tradeoffs

- $\uparrow$  block means  $\uparrow$  spatial locality usually
  - $\times$   $\uparrow$  block size =  $\uparrow$  miss penalty as we have to transfer more data
  - $\times$  Too little cache blocks in cache also bumps up miss rate (big blocks small cache)
- ~~any changes~~

	hit time	Miss penalty	Missrate
Size L $\rightarrow$ lf	inc.	Same	dec
Assoc L $\rightarrow$ lf	inc.	Same	dec
Blk size L $\rightarrow$ lf	Same	inc.	dec
Replacement policy FIFO $\rightarrow$ LRU	Same	Same	dec

Caches contd

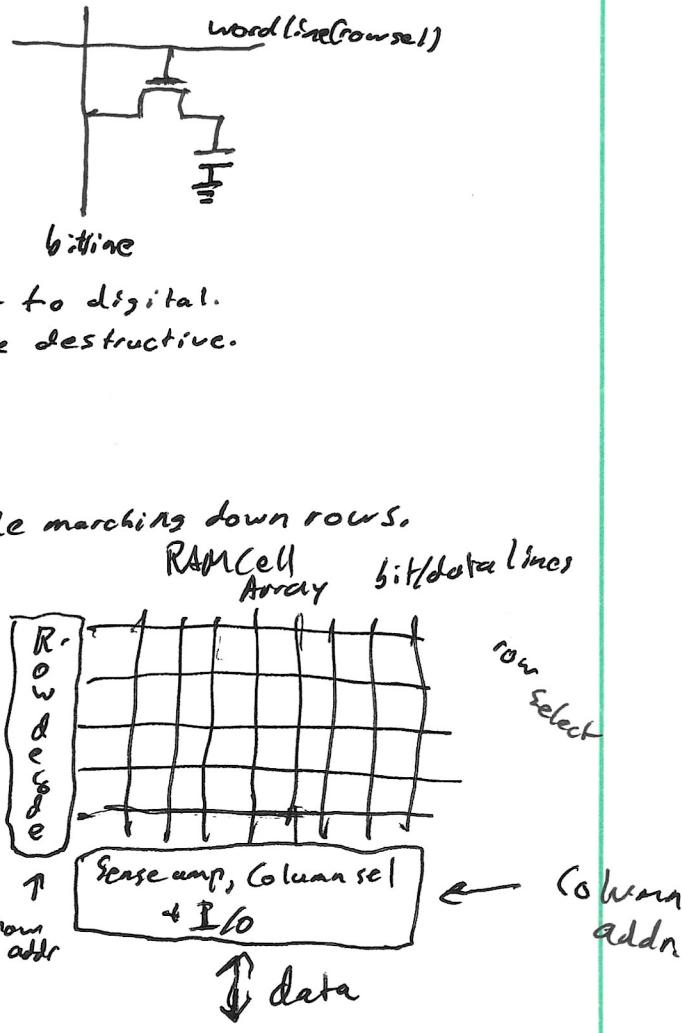
- Improving Cache Performance
  - $\downarrow$  hit time  $\rightarrow$  small example,  $\downarrow$  associativity
  - $\downarrow$  miss rate  $\rightarrow$  cache size, block size,  $\uparrow$  associativity
  - $\downarrow$  miss penalty  $\rightarrow$   $\downarrow$  blk size
- Need to optimize all cache factors and comprise with workloads in mind. (Simplicity often wins)

Cache Factors to optimize.	
- Cache size	
- Block size	
- associativity	
- replacement policy	
- write back vs write through	
- write allocate (bring block in cache on write miss)	

- Multi-level caches  $\downarrow$  miss penalty since miss penalty of L1 is just the AMAT of L2 instead of RAM latency.
- Split caches for instr + data can avoid structural hazards
  - instrs are 75% of mem. access.
- Modern CPU caches
  - 3 levels of SRAM cache
  - Low assoc. split L1: i + d cache
  - ~~high~~ higher assoc. in L2 + L3, 80% onchip transistors for caching.
- Memory Storage Elements
  - registers, latches - small, in processor, expensive to add
  - SRAM (caches) - medium, on chip close to processor, costly
  - DRAM (main mem.) - Large, off-chip, 50ns access time (2006 35-45ns for 512MB)
  - Disk etc. - Large, far from processor, slow(ms), (2006 20s for 200GB)
- Random Access - equal time is the same for all loc.
  - DRAM - Dynamic Random Access Mem.
    - $\times$  low density,  $\downarrow$  power, cheap + slow
    - $\times$  Dyn: needs to be refreshed regularly even if powered
  - SRAM - Static Random Access Mem.
    - $\times$  low density,  $\uparrow$  power,  $\uparrow$  cost, fast
    - $\times$  static: content lasts forever if powered.

Caches cont'd

- Not so RAM (Disk/CDROM)
  - Access time varies from loc. to loc. & time to time
- Sequential Access Tech - access time linear in loc. (tape)
- DRAM + SRAM (caches + main mem + reg) are not implemented as a bunch of mixed bits
  - Implemented as ~~arrays that~~ are read/written
  - DRAM is slow bc of 1 decode and 1 array, as well  
Long off-chip wires (mem. bus) between mem. + CPU.
- DRAM - Dynamic RAM/ main mem / RAM
  - dense, 1 transistor/cell
  - diff tech from CPU
  - forgets after awhile  $\sim 50\text{ns}$
  - impl. through ~~RAM~~ DIMMs (SIMMs (64b/32b transfers))
  - 1T DRAM cell = charge on cap.
  - Read
    - Precharge bitline to VDD
    - Select Row
    - Cell + bitline shares charge
      - cell cap charges, tiny Vchange on bitline
      - sense changes (analog)  $\xrightarrow{\text{bitline}}$  convert to digital.
      - write - restore val (reads are destructive.)
  - Write
    - Drive bitline
    - Select row.
  - Refresh
    - Do dummy read + restore while marching down rows.
- Classic DRAM organization.
  - each intersection is a 1T cell
  - 1 array = 1 wires so drain caps on each horz. wire + gate caps on each vert. wire
  - Adding pins is expensive
    - row + col addr time muxed (Same pins, latch vals.)
  - use 8 chips, 1 bit per chip to retrieve 1 byte.
    - can be expanded to larger values like 64b. access



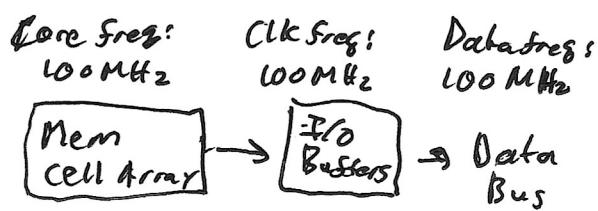
Caches contd

## • DRAM Optimizations

- Fast Page Mode - <sup>(FPM)</sup> consec. accesses to same row (row hits)
  - \* row Once & vary col. addrs. as row access is slow compared to columns (burst mode)
- EDO DRAM - extending data out - FPM + pipelining of col. accesses
- Synchronous DRAM - tied to sys clock, ↑ bus speed.
  - \* DRAM-DDR, DDR-2, DDR-3
- RDRAM - rambus - not standard.

## • DRAM organization - org./data transfer optimization.

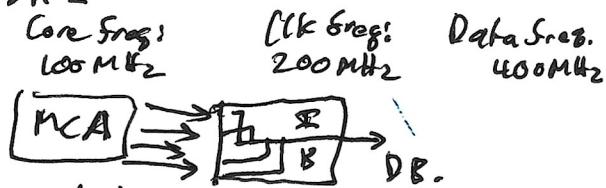
## - SDRAM



## - DDR1



## - DDR2



## • SRAM - static, no need for refresh.

- 6T implementation, on-chip, ~~off-chip~~ CMOS
- two inverters in a loop to remember WL access transistors.
- Slow bc of huge decode + array (smaller than DRAM) + long on-chip wires between caches + pipeline.
- 6T SRAM Cell

## X Write

- Drive bitlines ( $\overline{bit} = 1, \overline{\overline{bit}} = 0$ )

## X Select Row

## X Read (not destr.)

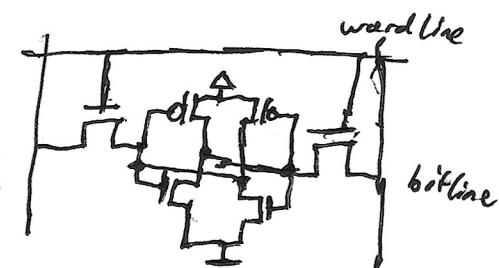
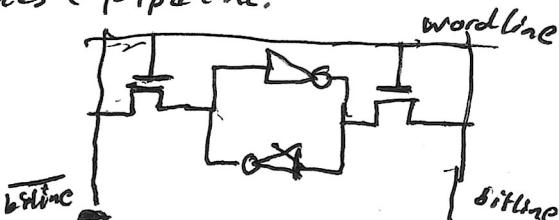
- Charge bitlines to VDD

## • Select Row

- Cell pulls bit or bit low

- Sense amp convert to digital

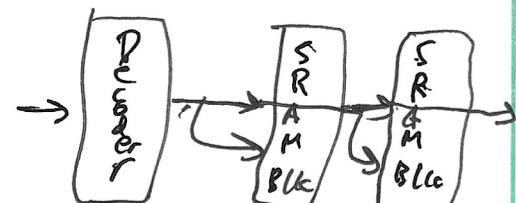
## - SRAM Organization



X Use similar cap policy to DRAM.

Caches cont'd

- SRAM Organization
  - Array's might be slightly different
  - Fully custom, b.c. no pins + on chip
  - wider access is better since first byte access is slow but consecutive bytes only small incremental cost. (wide access of whole block)
- Bank organization. - allows mass parallel access of blks
  - one huge array is ineff so SRAM+DRAM uses separate "banks" (arch. level)
  - banks are in chips + concat chips to get access width for DRAM
  - banks divided into sub-arrays for latency/area/power on the circuit level
    - SRAM uses <sup>L3</sup> Cache blocks in parallel for bandwidth
    - separate on-chip wires per bank, & access per bank, but still each access is slow (latency)
- Terms
  - Latency - total delay seen
    - \* Access time - time between reg. + word arrival.
    - \* Cycle time - time between reg.
  - Bandwidth - size of mem. transfer (T<sub>CO</sub> + Large Block Miss Penalty(LB))
  - Main mem = DRAM
    - \* address divided into 2 halves
      - \* Row access strobe (RAS) or column access strobe (CAS)
  - Cache = SRAM - no addr division
- Cache aware programming
  - Instr. Sequencing
    - \* Loop interchange: adjust nesting of loops for data access to match memory
    - \* Loop Fusion: combine 2 indep loops w/ same loop + some overlapping vars.
    - \* Tiling (Blocking): Improve temp. locality by accessing same blocks instead of accessing rows/cols.
  - Data ~~Layout~~ Layout
    - \* Merging arrays = ↑ sp. locality by joining arrays into compound elements
    - \* Nonlinear Array Layout: Map 2D array to lin. addr. space
    - \* Pointer-based Data Structures & node allocation



Multicore

- 3rd of 4 big ideas of modern computer systems.
- Why Multicore
  - other improvements  $\neq$  clk speed which means  $\downarrow$  power
    - $\times P = CV^2f$  ( $V$  linear to  $f$  so  $P \propto f^2$ )
  - lower voltage = leakage + reliability at some point.
  - $2x$  cores =  $2x$  power vs  $2x$  freq =  $8x$  power
- Parallel programming needed for multicore (not hidden from programmer)
  - way harder than seq. programming bc programmers have to track time, or which computation/who produces needs data when.
- Various Multicore methodologies
  - Shared I/O (Network) - use message passing (distributed systems)
  - Shared Memory - shared memory (primary focus)
  - Shared processors - single instr. multiple data/data parallel (SIMD)
- Parallel processes (threads) need to comm. (speeds cannot be assumed)
  - Message passing + shared mem. - programmer writes explicitly
  - SIMD is parallelized by compiler or written by programmer (programs)
    - $\times$  Works well if applicable.
- Shared Memory - all threads see one global memory (programmer view)
  - use loads/stores to access data
  - conceptually similar to uniprocessors + ease of programming (comm. is hard)
  - $\downarrow$  latency communicating small data items (often used by real programs)
  - hardware controlled sharing allows automatic data motion
- Message Passing - each process only sees its own private mem (programmer view)
  - uses sends/receives to access data
  - very simple hardware (almost nothing extra)
  - communication protocol is and must be explicit
  - shared mem. MP can emulate msg passing easily
  - has the highest programming burden, keeping track of comm. artifacts even w/ dynamic comm. patterns
  - have to carefully distr. data to processes
  - partition data or replicate (do repl. in software if not auto.  $\downarrow$  work)
  - coalesce small msgs into large msg to scatter to all priv. mems.
- Cannot assume thread relative speeds due to cache + branch preds.
  - Stalling all threads on miss.
  - Letting threads slip + slide relative to each other ~~allows~~ allows for useful work to be done on other threads during miss.

Multicore ~~Dated~~

- Shared Memory - hardware impl.
  - each thread has its own PC, registers, + stack (thread priv. vars are in reg + stack)
  - Rest of memory (heap) is shared (global vars)
  - Keep track of num of threads, threadids, + locks.
  - Use macros for creating threads, waiting for threads to finish, barriers to ensure thread synchronization, + lock + unlock to prevent race conditions from introducing issues.
- Message Passing - hardware impl.
  - all variables are private
  - Keep track of # of processes
  - use macros for creating processes, waiting for processes, send + received msgs.
- SIMD - hardware impl.
  - Shared mem/msg passing are SIMD
  - Instr on multiple data, & instr. overhead.
  - Deep, fast pipeline (ex. 50 stages)
  - Parallelism guaranteed by programmer/compiler, no need to check for hazards
  - Works great in some workloads (graphics/math核算 in AB) but not useful for all applications.
- Parallelization Strategies
  - Best parallel alg. isn't always the parallelization of the best seq. alg.
  - Think of partitioning tasks/data
  - Programming Models
    - x Pthreads - shared memory processors
    - x MPI (message passing interface) - clusters
    - x Open MP - shared memory - marking loops as parallel loops
    - x Thrust primitives - streaming, map reduce, parallel recur. (cycle())
      - customized to arch. like CUDA for nvidia GPUs.
- Typical multicore
  - Private L1, L2, L3 cache, shared banks of L3 cache
  - 64 bit mem. channel + multiple mem. controllers
  - mem + L3 connected to private cache w/ interconnect network.
- "Heavyweight Process" - diff thread of control (PC, reg, stack) + diff addr space
- "Lightweight Process" - threads - diff thread of control (context) + same addr space
- Can share regions across lt processes.

Multicore Contd

- Cache Coherence - important to maintain shared mem. Model
  - hard to track arbitrary sync among arbitrary threads in large distr. sys.
  - so all writes update all threads, to ensure to prior value isn't readable (write atomicity) without tracking sync.
  - value read must be val. from last write
    - x need to decide how to propagate writes + write serialization (single loc. write) or atomicity.
  - appearance of coherence not true coherency needed (constant updates) <sup>not acc.</sup>
  - coh. only provides write atomicity (no one can see old val.
    - x writes aren't instantaneous so these rules are followed
      - After writes, nobody can see old value
      - During writes, nobody can see new value
      - Later thread shouldn't see old val. while older thread sees new val.

## • Enforcing Coherence

- On writes, invalidate or update other copies
- Use snoopings in other caches to ensure other caches are updated and are not hitting on stale data
- All caches can see the bus + bus transactions to enforce coherence.
  - x cache misses go on bus to snoop in all other caches (looks at tag)
  - x cache controller generates block states on processor/snoop events + generates bus actions.

## • Shared Mem. Models

- Shared Cache - processors share a cache which is where all data goes through.
- Bus-based shared mem. - covered here
- Dance hall - similar to bus based but uses an interconnected network + pt to pt connections.
- Distributed Shared Mem. (DSM) - memory local to each processor w/ an interconnect network.

## • Snoop Design Choices

- set of states
- state transition diagram
- actions
- write-through vs write-back
- update vs invalidate
- multiple-reader, single-writer.

]} = Protocol

]} = Basic Choices

Multicore

- 3-State write-back invalidation protocol
  - MSI protocol
    - Modified - stale mem, + one cache has latest valid copy
    - Shared - 1+ caches have valid copy + no stale mem.
    - Invalid
  - invalidate all copies <sup>Other</sup> before M state.
  - use bus tracking to enforce
  - MSI States
 

The diagram illustrates the transitions between three MSI states: M (Modified), S (Shared), and I (Invalid). Transitions are labeled with bus requests and writes:

    - M to S:** Bus Rd/Bus WB → S
    - S to M:** Bus Rd/Bus RD → M
    - S to I:** Bus RD/Bus RD → I
    - I to S:** Bus RD/Bus RD → S
    - I to M:** Bus RD/Bus RD → M
    - M to I:** Pr Wr/Bus Rd → I
    - I to M:** Pr Wr/Bus Rd → M
- (is invalid if no tag match)
- Pr Rd/-  
Pr Wr/-  
Pr Rd/-
- Pr Wr/Bus Rd
- Pr Rd - curr. processor Rd
- Pr Wr - curr. Processor Wr.
- Bus Rd - some processor rd.
- Bus Rdx - some processor rd. lwn.
- Bus WB - some processor writing back data
- A/VB = if A, B = A → B
- Three edge cases
  - × S → M (Pr Wr/Bus Rd) - we have block but mem. Still provides blk.
    - need Bus Rd for invalidation,
    - need to read from mem on Bus Rd for I → M.
      - can use separate Rds for I → M / S → M, but it's more complicated
      - On Cache side, act as cache miss (Bus Rd) or Pr Wr to S.
  - × M → I (Bus Rd/Bus WB) - bus writeback even though other cache blocks bigger than word, <sup>Processor is overwriting</sup> so need to wb entire blk.
    - Mem sees Bus WB + doesn't know if M → S, even though mem will be stale at the end
    - can make sep wr for M → S / M → I, but more complicated
  - × S → S (Bus Rd/-) - can provide block to read.
    - multiple Ss might exist, leading to need of arbitration + complexity.

Multicore Coher.

- Details
  - Each cache block has its own state + multiple caches share state (if same blk)
  - ↳ Separate cache controller (per cache), bus controller, + mem. controller for clean & easy impl.
  - BusRd + BusRx have responses (blk), BusWB has no reply
  - Op. sequence:
    1. core req. (dst)
    2. cache hit/miss (cache controller)
    3. bus req. (bus controller)
    4. grant + bus transaction (BusAd/BusRdx on bus)
    5. snoop into other caches (each cache controller)
    6. check MUC state of all caches in parallel (each cache controller)
    7. perform state trans. + update state (each cache controller)
    8. access mem. if needed (mem. controller)
    9. either memory or a cache place block (BusWB) on bus
    10. requestor cache puts blk in cache, updates state, complete dst (co.)
    11. data returned to core if needed.
- Bus Controller - maintains coherence + broken into multicycle to maintain fast cycle.
  - has separate state machine from cache controller.
  - 5 bus cycles
    1. Bus req. (bus controller)
    2. Arbitrate + bus grant (bus controller)
    3. Snoop + state transit (all cache controllers)
    4. resp. data (wait many cycles if from mem) (1 cycle/memory)
    5. req. compl. + release bus (1 cache controller)
  - Must latch important signals in the cycles to prevent long delays
    1. latch busreq. (addr, trans, ...)
    2. latch bus grant (addr, trans, ...)
    3. latch snoop resp. signals
    4. latch resp data on bus.
    5. latch resp data in cache.
  - Atomic bus - not necce. (can pipeline) but is way easier.
  - Cache mem. access chart will be modified to incl. bus controller
  - use a duplicate tag array for snoops for no structural hazards between pr + ~~bus~~ cachebus snoop.
    - ↳ must update tag array + cache together, but it's ok to tag updates only on eviction/write.
    - & stretch state update till both tags can be updated.

Multicore Cont'd

- Coherence makes writes appear atomic
  - hits guaranteed correct data
  - services one miss at a time in order received to ensure two parallel write misses
  - Completes servicing entirety of one miss before moving on bottlenecked @ bus arbiter.
  - real systems do not use atomic bus (slow)
    - x uses split transaction bus (divide req + resp or pipelining)
      - write atomicity is complex - used in small systems
    - x uses no bus in distributed shared mem. (DSM)
      - write atomicity still complex - used in modern processors
- Invalidate vs Update
  - Invalidate - one write before many reads
  - Update - many writes before one read.
  - no clear winner, but Update wr. atomicity is harder in DSM systems, so invalidate used in practice, not update
- Synchronization ensures data races don't occur.
  - use lock + mutex to ensure one thread only can run "critical" section at a time.
  - barrier ensures thread sync before continuing
  - event notification as cross thread signaling
- Atomic R/W ops ( $r \neq \text{lock\_var}$ ) - op has to be atomic & return prev lockval( $r$ ) in all cases
  - Test + Set ( $r, x$ ) ( $r \neq m[2], m[2] = 1$ )  $r$  stores lock & checks if  $= 1$  (already locked)
  - Swap ( $r, x$ ) - swaps  $r = m[2]$ ,  $r$  set prev. lockval.
  - Compare + Swap ( $r_1, r_2, x$ ) ( $\text{if } r_1 = m[2] \text{ swap } r_2, m[2]$ )
    - x  $r_2$  get prev lockval
  - Fetch + Op ( $r, r, op$ ) ( $r \neq m[2], m[2] = op(m[2])$ ) -  $r$  gets prev lockval.
- Lock implementation
  - basic! Lock:  $\text{while}(r \neq \text{lock\_var}) == 1;$   
Unlocks:  $\text{lock\_var} = 0;$
  - enforcing & doing atomic rmw in hardware is hard, so it's easier to detect non-atomicity & retry till it works
    - x can detect by writing when locking (invalidation of other copies)
  - use lr-sc instead of lwsr for atomicity
    - x lr - load-reserved - stores reservation set addrs next to cache
    - x sc - store cond. - stores data & invalidates other copies if rs. not invalidated
      - invalidate then check for consistency. (invalidate is too expensive)
    - x one res. set per core.
      - x no store of sc if rs. invalidates, update res accordingly (R)
  - atomicity req. for a lot of OSs

Multicore cont'd

- All sync. stuff is block addr. granularity
  - coherence (now is blk gran.)
  - v.s. addr. - two locks in same block is bad perf - false sharing  
but does work
  - If Id est done to lock it's a software bug.
- Performance
  - T+T is nasty for high contention locks due to constant writes  
x writing to all for non-holders
  - a lot of cache misses due to large # of invalidation
  - can optimize by ↓ T+T exec. + operate + reduce bus traffic
    - x T+T w/ backoff + insert exp higher delay between T+T
    - not fair as later cond. might get earlier
    - x T+T+T keep testing + only T+T when lock is released
      - no cache misses while waiting but still unfair + contention
- Synchronization
  - Deadlocks - hard to break out of w/o external intervention
  - Races - wrong/no lock used leading to racing for data.
  - Timing-dependent - super rare interesting and hard to debug
  - Lock composition - hard to divide program into sep. locks.
  - Transactional Mem - not currently ~~concurrent~~ + uses blocks of ops in one atomic op.
- Coherence vs Consistency - coh. is for one mem. loc., cons. is for all loc.s
  - coherent if ops performed by processor appear in prog. order
  - coherent dictates what consistency dictates when (ordering)
  - reg alloc. in compilation can ruin SC, use volatile to remove reg assignments
  - consistency model is deal between sys. + programmer  
that tells programmers what behavior to expect.
- Sequential Consistency Model (SC) - programmers think like this
  - Accesses from one thread in prog. order & writes done atomically.  
x no rel. speed assumed + each access complete before next starts
  - Sys processors issue mem ops in prog. order.
  - Processors take turns randomly to make random # of mem accesses.  
& gives illusion of sequential thread handling.
  - don't preserve strict mem order for of performance, in software issue w/ parallel prog, other sols could exist.

I/O

- Needed for users, computer, and non-vol. mem. communication.
- Total time = CPU + I/O - overlap
  - Transaction Processing does many indep. small I/O ops.
- Performance considerations of I/O depends on applications
- Disks - trends have ↓\$/MB, ↓ disk diameter, & seek time, parallel → serial
  - service time
    - seek time - ~10-20 ms (better w/ locality)
    - rotation time (sector) -  $\frac{1}{2}$  rot. on average - 4ms - 8ms (3600-7200 rpm)
      - ~~4ms~~
    - transfer (sector) time - Pg size transfers ~4-8KB (several MB/s/s)
  - access time = seek time + controller overhead + service time.
- I/O reliability/dependability/availability
  - MTTF - mean time to failure - reliability/availability metric
  - MTTR - mean time to repair.
  - MTBF - mean time between failure -  $MTBF = MTTF + MTTR$ 
    - availability metric.
  - Bit flips - corrected w/ ECC bits for single/double error detection.
    - done by disk/mem/cache controller
  - Disk errors - easily fixed w/ parity
  - Using banks of data helps error fixing
    - group entire block in one bank & spread it to different ~~blocks~~ banks.
    - getting first byte is slow so this is better than spreading blocks.
  - Redundant Array of Inexpensive Disks (RAID)
    - provides reliability through redundancy & ~~parallelism~~ performance
    - RAID 0 - no redundancy through parallelism.
    - RAID 1 - mirrored for 100% overhead & redundancy
    - RAID 3 - bit interleaving parity - small overhead.
    - RAID 4 - block interleaving parity - small overhead & writes
    - RAID 5 - distr. blockwise parity - small overhead, writes, w/o bottleneck.
- Connecting I/O to CPU w/ busses. - needed for connecting w/ ~~any~~ I/O.
  - Adv.: versatility, low cost, can manage complexity
  - Disad.: com. bottleneck, speed limited by len & # of devices on bus.
    - need to sup. varying transfer rates & latencies.
  - Processor - Mem bus (design specific) → short, high speed, maximizes data transfers from mem.
  - I/O bus (industry standard) - lengthy & slower, matches variety of I/O
  - Backplane bus (standard/not) - interconnection structure for all components, connects to another bus.

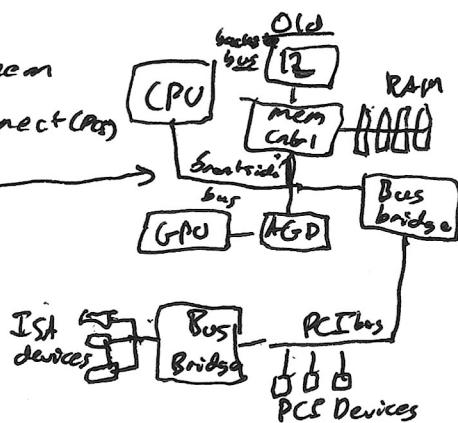
I/O cont'd

## • Bus Architectures

- one bus (backplane) - all devices on bus
  - + simple + low cost
  - slow + can cause bottlenecks
- two bus - uses bus adaptors + I/O buses that interfaces w/ proc-mem bus.
- three bus - uses a bus adaptor + backplane bus to connect to further bus adaptors + I/O buses to connect to proc-mem bus.
  - x reduces loading on proc. bus.

## • Real World Architecture

- Really old - dual indep bus between proc. + mem
- Older - Peripheral Component Interconnect (PCI)
  - x 66 MHz sync, 32bit data, multi master
  - x AGP used for high speed connection to FSB for GPU
  - x parallel
- Modern - PCI-e (express)
  - x serial transmission + faster w/ pt to pt.



## • Bus definition

- Transaction Protocol
- Timings + Signaling Specification
- Bundles of wires
- Electrical Specification
- Physical / Mechanical Characteristics (connectors)

## • Synchronous Bus

- + little logic + runs fast
- must run at same clk for all devices + must not be long if fast to avoid clock skew
- often uses CLK, bus req, bus grant, r/w addr, wait, + data signals

## • Async Bus

- + no clock skew
- can be lengthened + accommodate a variety of devices
- requires handshaking protocol.
- often uses readreq, data, ack, data ready signals

## • Flash memory - non-volatile mem, faster than disk, slower than DRAM

- electronic storage - more reliable, faster & costlier than disk.
- writes slow bc of wearout (1nJ/cell)
- r/w in blocks.

## • Memory-mapped I/O

- maps I/O devices to VA w/ address + data buffers / queues
- I/O is shared w/ only access to OS which sets up VT + PA mappings.
  - x unreachable for obs. reasons

## • I/O devices need to notify OS of when op overflows.

- Polling - period checks to status reg., simple + only processor controlled
  - x but overhead = TCRU time (ineff.)
- I/O interrupts - only halts during transfer, but T interrupt overhead hardware
  - x sync
  - x need to comm. status + device
  - x can be handled when convenient, but need to respect interrupt priorities.

I/O Control

- Direct Memory Access (DMA) - allows for data transfer from mem.  $\xrightarrow{\text{CPU}}$   
 - acts as master on bus.

- uses phys. mem addr. (might be translated by OS)

- one page per transfer (pinches to prevent pg faults)

- I/O processing

1. Issue instr to I/O processor (IOP)

2. I/O looks in mem. for commands

3. Devices to/from mem transfers are controlled by IOP. (steals mem cycles)

- I/O setup

- before cache slows CPU, so can't do

- after cache requires coherence. (applies to DMA too)

- CPU involvement - low to high involvement

- x Polling - does everything: detection, I/O, copying, mem.

- x Interrupts - on notification, does the work

- x DMA - CPU describes work, DMA does it.

- x IOP - CPU tells IOP where to find instrs.

- OS responsibilities

- x OS interfaces between proc + I/O for protection.

- x Provides abstraction, equitable access, handles interrupts, & schedules accesses.

- Multiprogramming:

1. Proc. invokes syscalls

2. OS checks prot & runs device drivers

3. Suspends curr. process & switches process

4. I/O Fielded by OS

5. OS completes I/O & wakes up suspended process

6. Next process that is ready to run.

- Video processing

- Frame buffer - stores frame info (color b/pixel)  $\xrightarrow{\text{FB}} \xrightarrow{\text{Render@FR}} \xrightarrow{\text{VGA}} \xrightarrow{\text{main mem. bus}} \xrightarrow{\text{Monitor}}$  #1

- Setup #1 is dumb + slows bus

- using VRAM sep. to main mem

- x eliminates structural hazard (#2)

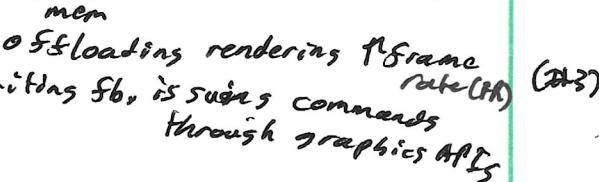
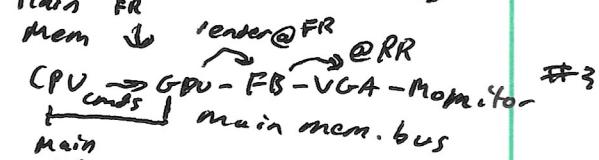
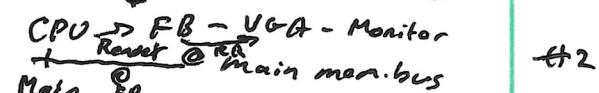
- x dual ported w/ VRAM

- x serial port for RasterScan w/ shift reg.

- using independent video accelerator for offloading rendering frame

- massive reduction + bandwidth (GPU)

- x APIs like OpenGL & Direct3D



- Other accelerators

- Per application (ex. GPU for 3D modeling)

- Broader applications (ex. using GPGPU ( SIMD parallelization) for computation)

## (SMT)

SMT

- Simultaneous Multithreading / Hyperthreading (Hardware MT)
  - stalls happen in pipelines even in out-of-order (OoO) processors (cache miss)
  - peak 6-8 IPC, actual 1-2 IPC (25-35% eff)
  - overlap thread execution during stalls
- SMT gives instruction level parallelism
  - run 2-4 threads (or processes) simultaneously on same pipe
  - not context switching in OS but thread switching in software.
  - thread level parallelism - overlapping mem. references to improve throughput.
    - x latency of each thread is same or slower
    - x hardware utilization is better as we overlap cache miss w/ other threads CPU work.
    - 2-4 times more requests to mem but w/ DDR versions, bw will be enough.
      - x does cause capacity pressure on caches.
- SMT implementation
  - need 1 hardware context per thread (multiple instr in pipe at once)
  - some hw will be duplicated, some will be shared.
    - duplicate:
      - x PC, reg file, wr. buffer, PTBR, reservation set, PBD (for correctness)
      - x Br prediction table (for performance)
    - Shared:
      - x fetch, decode, Ex units, mem, larger caches, TLBs
  - OS views ~~the~~ contexts as a virtual CPU
    - x assigns processes to each core for all virtual CPUs (threads)
  - SMT allows sharing of most of the pipe for hardware utilization.
    - x multicore, each core is usually SMT capable.
    - key decision is where/which context to fetch
      - x rest of pipe doesn't matter other than, reg r/w, write buffer, PTBR, & fetches from non-stalled context + squashes pipe when stalled, reset, PBD
- SMT will be correct memory/caches/TLB as ~~all~~ processes will have diff phys addrs. thanks to virt. mem.
  - diff. threads of same process go to same phys. addrs.

## • Latency

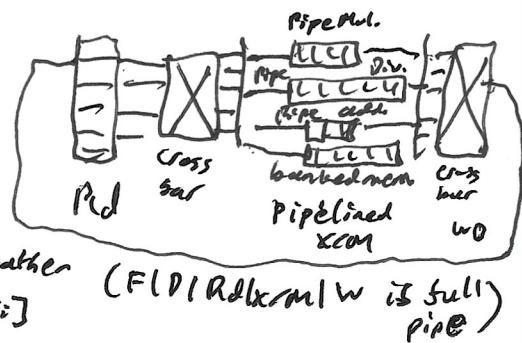
- (2-4/  
12-16 cycles)
- L1 & L2 hit - overlapped by OOO in 1 thread

(40-60/  
~300 cycles)

  - L3 hit + mem hit - overlapped by OOO in 1 thread + SMT + Multicore (across multiple millions of cycles)
    - Mem miss / page fault - overlapped by OOO issue in 1 thread + SMT + Multicore + OS multitasking (across multiple threads/processes)

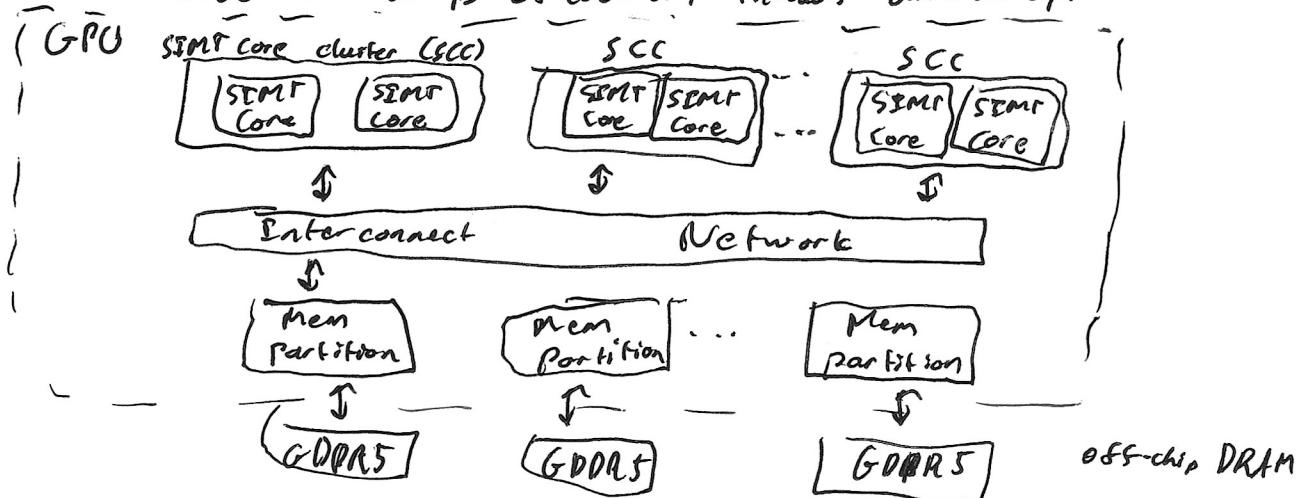
Vector

- Used for parallel regular code like matrix calc in scientific domains
  - CPU made for gen purpose ugly code and thus is complex.
  - CPU complexity makes a separate vector handling structure more effective.
- Vector has highly parallel regular code and thus wants deep or wide pipe
  - deep = fast clk, wide = #parallelism.
  - no data or inner loop hazards <sup>control</sup> so can avoid that h/w
  - multiple ops per instr → statically-sched. in-order wide/deep pipe.
  - uses data level parallelism (DLP) through single instr. multiple data (SIMD)
- High level code → vectors by several processes,
  - Compiler w/ automatic vectorization
  - Programmers w/ hand tuned assembly
    - High level lang. like CUDA
      - hard to prog. close to assembly level
      - huge perf drops if prog. doesn't match h/w
      - done to optimize tasks really well (like AI)
- For handling vectors wider than SIMD through
  - Pipeline ops in one deep pipe / one wide pipe (multiple func units)
  - use vector reg to store values after retrieving from mem. subsystem
    - also writes back to mem. subsystem at the end (intermediate stored in reg)
- Vector Mem. system
  - Load/Store from mem into vregs (operates on contiguous addr.)
    - assumes a unit/reg stride length according to vreg/data size
      - Ex.  $M[r1 + i \cdot r2]$
  - Has lot of spatial locality but not temp. locality. (contiguous streaming access)
  - compute so compute per byte
  - System great for high bar DRAM that hides latency w/ prefetching.
    - heavily banked for bw
    - uses to hold prefetched data.
- Strip mining - used if data > vec reg size
  - divides loop into smaller ops & vectors.
  - can programatically set VLEN for ops.
- Indirect Access - use vector indirect scatter/gather
  - custom instr. (ex.  $M[r1 + v2[i]] \rightarrow v1[i]$ )
  - (FLD/RDLM/W is full)  
PIPE
- Branches - use vector masks
  - hides vector based on cond. (ex.  $\text{mask}[i] = 1 \text{ if } v1[i] < 50$ )
- Vector extensions created for processors
  - ~~Euclid~~
  - commodity vector processors - embedded, DSP, GPU, etc.



GPU

- Started for graphics but GPUs now run all kinds of SIMD workloads like ML, Vision, etc.
  - ↑ perf per watt + much more (10x) energy eff than multicore per op.
  - runs a SIMD microarch (diff from normal vector pipes)
  - SIMD processing w/o SIMD instructions
    - execute scalar threads in SIMD exec or SIMT
      - more prof than vectors, but only works if SIMD-amenable
      - software should be SIMD friendly for good perf.
      - hw width hidden from software SIMD
      - hw group = warp in NVIDIA.
- GPUs - graphical processing unit.
  - accelerator for raster based graphics (OpenGL, DirectX)
  - programmable commodity hardware w/ 1000s of ALUs and 10s of 1000s
  - Exploits TLP + DLP (fundamentally vector processor)
  - uses multithreading to tolerate mem. latency ~~concurrent~~ threads concurred
  - Software written for single thread
  - ↑ programmable than vectors (using threads) w/ CUDA
  - ~~SIMT~~ programming - single instr. multi thread.
    - fine grain threading unlike CPU SIMD.
- GPU used for compute bc it uses way more % of silicon for compute compared to CPUs
  - also uses significantly less energy per op.
  - achieved by having minimal cache due to streaming access
    - x cache recently for ML workloads.
  - run code w/ warps of lock step threads divided up.



GPU Cores

- GPU pipelines - wide & deep unlike fast deep vector pipes
  - wider for slower clk + energy eff.
  - super simple efficient pipelines w/ 32 lanes in lockstep
    - & each lane has regd, exe, mem, wb stages
  - no bypass/bypass by hiding latency along w/ mem latency w/ multithreading
  - inorder issue, out-of-order complete.
- Register file - do one wide read for all lockstep threads
  - 4 banks of res access is enough bc of this
  - read one res at time when overlap of access
  - only avg case matters as worst case is hidden by MT.
- Branches - cannot use vector masks to keep programmable
  - lock-step lanes assigned if taken/ is not taken
  - run is taken first then is not until convergence (found by compiler)
  - hardware stack for nested if/else - tracks divergence & convergence
  - too many data dep. branches = ↑ divergence & perf
    - & graphics often doesn't have data dep branch, so it works.
- GPUs use small L1 cache due to not much reuse b/c streaming workloads
  - use coalescing of mem access instead of 32 ports to only make 1 block access. (otherwise perf ↓)
- NVIDIA calls each 32-wide pipe as a streaming multiprocessor (SM)
  - each SM has L1 + software managed scratchpad called shared mem.
  - L2 shared by all SMs
  - Mem = global mem.
  - can be 32wide 32b, 128wide 8-bit, etc.
  - each chip might have 25-80 SMs
  - lot of access is predictable but not all (texture caches)
- Unpredictable access exposes latency addressed by MT.
  - GPUs have gazillion threads, ex. 32wide pipe w/ 50way MT.
  - Run any ready context every cycle, schedule w/ scoreboard.
  - High MT possible in GPU due to extremely high parallelism,
    - & in programs + data
- GPU mem. - ↑ bw & ~~low~~ MT
  - 50 contexts each 32 wide per SM
  - 32wide usually coalesced access to one blk + many diff contexts on one SM often access same DRAM row (mechanical).
  - ↓ locality + ↑ streaming → more banking for bw.
  - can tolerate almost any latency.
  - no coherence / freq. sharing, load+comp. flush sw, low VM supp. in GPU vrs.
- No coherency sharing, or fine grain comm. (not needed)
  - Comm. over PCIe or in integrated GPU.
- Tensor Cores added to GPU for specialized matrix mul.
  - often done in ML + runs way, way faster than SIMD
  - lots of SMs + parallel mul. w/ systolic arrays.