

- Exceptions** → Exception: unusual internal event caused by prog. during exec → interrupt: external
 → trap: forced transfer of control caused by exc. or interrupt → not all excp. cause traps
 → Asynch interv: I/O asserts an prioritized interrupt req lrcs, choose to process interv, stops at inst_i, completing until J_{i-1} (precise interrupt); & EPC to whichever didn't U/commit
 → Trap handler: save EPC before enabling interv, ERET (enable + restart) Sync trap: particular instr, usually restart
 S-stage: hold exception flags until commit point (M stage, Speculation burst)
 ⚡ only wr at commit, else throwing & flush & launch handler ⚡ Bypassing allows uncommitted instr use of results by following
- | | | | |
|-------------------------|--|-------------------------|------------------------|
| F F D X M (W) committed | raise → wait for all to get to FD- start | Imprecise: FD X M1 M2 W | Reduce latency: |
| F D X (W) NOT committed | X cache miss | FD X M1 M2 W | unroll, kill precisely |

Latency: t taken for single op to start to finish Bandwidth: rate at which ops can be performed Occupancy: t during which unit is blocked in op (struct. max.)

Memory DRAM (optimized for density) SRAM (optimized for speed), not destructive CPU ↔ SRAM ↔ DRAM ⚡ Temporal: for self, ref. again ⚡ Spatial: near locations 3 predictable forms

\$ Replacement ⚡ Random (Fast) ⚡ LRU (Complex) ⚡ FIFO (high assoc) ⚡ NMRU (FIFO except most recent)
 \$ Policies Cache hit ⚡ write-thru (\$+ncm) ⚡ write back (\$only, eviction spawns wr) ⚡ write-through wrthru+wrallc
 Cache miss ⚡ No-write alloc - wr to mem ⚡ write alloc - fetch into \$ (only to rem if evicted) ⚡ wr back + wralloc

C's: larger lrc size; ↓ compaction + conflict Larger \$ size: ↓ capacity/conflict ↑ hit% Higher assoc: ↓ conflict ↑ hit%

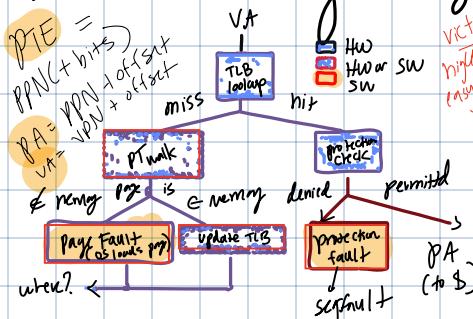
Write buffer: hold wr before \$, reduce read MP ⚡ reduce wr penalties (w/o) ⚡ cost of evictions Misses per str

⚡ Local missr = miss in \$ / access to \$ Global missr = miss in \$ / CPU num access (NP) = misses / total access

Prefetching Accuracy (useful/predicted) coverage (useful / total vng. access) Time (on-tvc / total perf.)

Algos: Next-lirc (next N lrcs after comp. miss) Strided (after say N w/ dist D → $\frac{N}{D}$) $b+2N$

Virtual Memory Page: fixed size, internal frag, efficient ⚡ Segment: variable, ext.



Ex 4 GiB VA, 4 KiB pages, 4 byte PTEs $\rightarrow 2^{VPN} \times \text{page size} = PA$
 ↳ Page offset = $\log_2(\text{page size in bytes}) = \log_2(4096) = 12$ size addressable
 ↳ VPN bits = Size of Mem addr - page offset bits = $32 - 12 = 20$
 ↳ # virtual pages = $2^{VPN} = 2^{20} \text{ PTEs} \star \text{VPN} = 2^{VA/\text{page size}}$
 ↳ 1 page mapped to PA → 1 valid PTE
 ↳ linear page table $\rightarrow 2^{20} (\text{PTEs}) \times 4B(\text{size}) = 4\text{MiB PT}$ L1 vs.
 ↳ 2 level: 1 page mapped to PA → 2 valid PTEs $\rightarrow 2^{10} \text{ PTEs} \times 4 \rightarrow 2^{12}$

VA: 0x58 → 0x5 = 0101 → 01 = index 1 (6 bit works → 0x00+8) → go to 0x040 → apply 12/3 → PA

Final PA: Content + offset Ex 0x17 content → 0x178 final ⚡ Bring in pages for faults

AMAT: HT: incr if crit path↑ MP: increases w/ Block size / outer num, MR: decr. w/ bigger tag

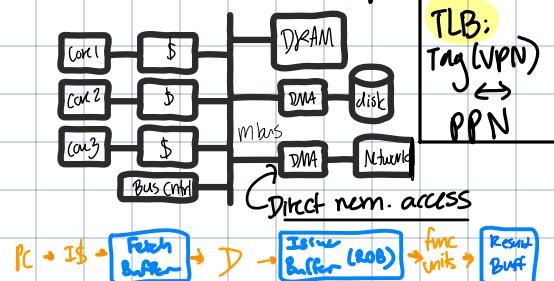
Ex Matrix: Bits 0 to 11 off set + index $\rightarrow 2^{12}$ VS. + 4096 bytes b/w col $\rightarrow 2^{12}$ index

TLB: \$ for PTEs, fully assoc ⚡ TLB miss ≠ Page Fault ⚡ TLB lookup = tag check

TLB entry: valid/dirty bits, tag, PPN \rightarrow PPN + offset = PA, use VPN for tag

VINT: cache index & cache tag derived from VA VIPT: VA for index PA for tag → parallelism

Digits	2 ²	4 2 ⁸	256	Ki 10
0000 0 0101 5 1010 A	2 ³	8 2 ⁹	512	Mi 20
0001 1 0110 6 1011 B	2 ⁴	16 2 ¹⁰	1024	Gi 30
0010 2 0111 7 1100 C	2 ⁵	32 2 ¹¹	2048	Ti 40
0011 3 1000 8 1101 D	2 ⁶	64 2 ¹²	4096	Pi 50
0100 4 1001 9 1110 E	2 ⁷	128 2 ¹³	8192	Ei 60



TLB: Tag (VPN) \leftrightarrow PPN

RISC instrs

lw rd, rs1	bge = bgeu \geq
sw rs1 rs2	bne \neq Slt <
mv rd, rs1	blt < Slti < im
bge	SLtu < u
bltu	SLtu < SLtuVim

SubRoutine Return Stack / Return Addr stack (RAS): store retr. addr as called, more acc than BTB
 (x) $f(a) \rightarrow f(b); f(b) \rightarrow f(c); \rightarrow [8FCC] \rightarrow$ ① Lookup BTB for jmp ② push ret. addr ③ return ④ pop stack ⑤ addr

InO: no instr issued after Br can WB before branch resolves **VS** DoO can complete DoO but cannot commit

VLIW (V. long instr word) force compiler to schedule so that HW doesn't have any arb interlocks

VLIW = N bundles, each is a type specific slot ④ 2 muls *** Guarantee intra-instr. parallelism**

Loop unrolling: mul iter in one big core, faster stride, a lot of overhead

SW pipelining: same # iter, only 1 time overhead, modify instr to hide latency **AT**

*** Write out + two iterations:** 1st enters, 2nd starts in loop, then last is readin/o loop

→ Best to introduce if MEM in loop: is not fully utilized (spurious wait)

[VLIW Issues: obj. code compatibility, wasteful size, unpredictable latency is dangerous to it br. \$
 ↳ Compilation techniques: static scheduling (fixed/no prediction), SW pip, unroll, trace scheduling

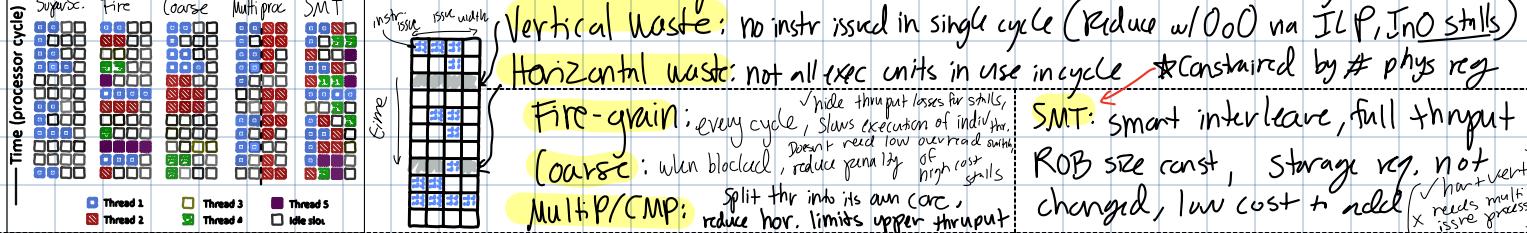
Predicated execution: Mispredicted br limit ILP → eliminate hard-predict br via "if-else encoded inb/instr"

Instr conditionally execute; avoid var. len latency in IF cond → do this else → NOP

Trace Scheduling / Trace: path of blocks for most taken br path, schedule whole trace at once

Multithreading (TLP) Incr. Thruput (more instr sent) worsen latency (competing concurrency for resources)

*** Doesn't benefit from br. resolution latency (can hide); arch. identical to multicore system to ISA**



i-count policy: for SMT, fetch thread that has least instr in flight (if stalling/in flight, will get stuck)

Sharing Resources vs. Dup: (Duplicated) PC, Rename, Fetch unit, phys Regfile, window, units, ROB

loop: 1d ₃ ss, A(S1) +1	Fire grain	Data-Dep	loop: addi a1 a1 -0x1 3NH
1d ₃ ss, B(S1) N+1	① Find max	① Ellipse	1w n0 o(a0) +1
1d ₅ ss, Y(S1) 2NH	② steady max	2 steady max = group	cycle breg a1 loop 2N+1
2addi s8, s8, -2 0NH	③ Find N	3 ceil (so-4) +1	Fire grain: +1 & 4N+1 same time
fmul st, s8, s8 3NH			Data Dep: ceil (so-3) +1
fadd ss, st, s8 1NH			since circles round!
sd ss, S(S1) 2NH			
2bnez s8, loop 0NH			

for (unsigned int i=0, i < N, i++) {

if (a[i > idx[i]] < 0) {

b[(dilation * i)] = scale * a[idx[i]];

Note: a[0]=n, ra=a[1], xb=a[2], y[idx]=a[3], sidx=a[4], as[5]

slli as, as[2] dilate to bytes

vsetvli fo, a0, e32, m1, t, m1

vle 32.v v1 (a[3]) load idx[i]

vsll.vi v1 v1 2 convert to byte offsets

vtxe 32.v v2 (a[1]) v1 load a[Idx(i)]

vmsl4.vx v0 v2 x0 set mask (idx[i]>0)

vmul.vx v2 v2 a4 v0.t Scale * a[Idx(i)]

vsse 32.v v2 (a[2]) a5 v0.t Store b[dilation * i]

sub w0 no fo decrement n

mul f1 fo as v1 * dilation in bytes

slli fo fo 2 v1 to bytes

add a3 a3 fo bump ptr to idx

add a2 a2 f1 bump ptr to b

bnez a0, loop

and: ref

VLMAX = $\min(MUL, SEW)$

LMUL = $\min(MUL, SEW)$

VA PA translated addr vs. using VA to address data in \$ → Mul VA + single PA (Aliasing: sim wr/reads)

VIPT: offset bits = virtual index → P tag & check *** Page offset = offset + index bits, else L2 to check**

Vector ISA

SEW=8 LMUL=2

vsetvli fo, a0, e8, m2, ta, ma (config vlen & vtype)

vise 32.v vsd (rs1) rs2, (rs1) (32 bit striped load)

vmul.vt v2, v2, x0 (set if less than (signed), v0.t = mask)

vtxe 32.v v2 (a1) v1 (load a[Idx(i)]) use elem of V1 to index into V2

vmul.vt v2 v2 v2 v0 v0.t (scale + a[Idx(i)]) IE cond. wise

vsse 32.v v2 (a2) as v0.t (store into b[dilation * i])

Vectors (DLP)

ISA separation from math

Compact, Scalable, Expressive

Packed SIMD: fixed vector len,

no shaded ls, st, must

aligned instr set, scatter/gath, align

key, sup. dispatch to store busy, loop unrolling ↑ reg pressure

Vectors (Cont.) Vector len register = # valid op⁺ MAX VL / VLMA ×

Standard elem width (SEW): default width of each vec elem

Len Multidiv (LMUL): int power of 2 grouping vectors $\times = 2^{vn \times v \times n+1}$

VL = # cores → 32 elem vector, 8 lanes → 4 vectors issued per instr

Mem System: (Banking) split address out, go to sep. bank w/ high probability

Chaining: bypass vals early; no chain if two R's but separate func units

Strip-mining: finite V size → split across mul vs; best when $ik \ll n > MaxVL$

Vector reduction: BinTree reduction, recursive aggregation

Sum += A[i] → sum = sum + A[i] dupmul A[i-1] + A[i-2] ...

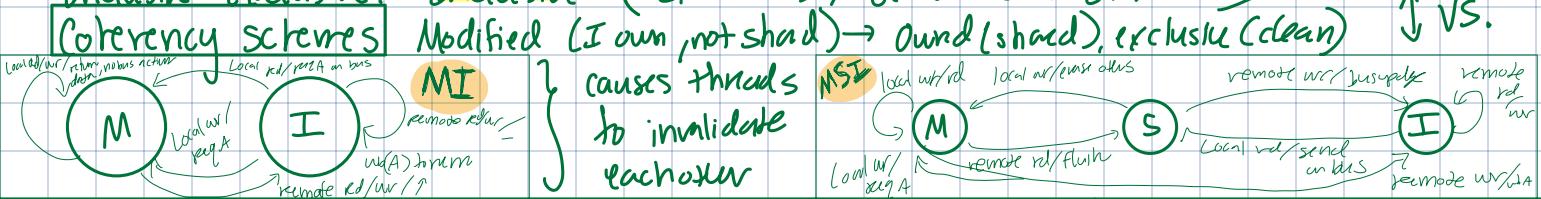
SIMD (vector) SIMD (thread) SISD (000) GPU: SJMT (multi+thr) = DLP × TLP

Thread blocks: range over iteration space Warp: sub blocks not visible to user, microarch usually

Final Cheat Sheet

\$ coherence: coherence: legal values a single mem adder can return ~~if not transient~~ if not transient
write policies: wr-back: only \$, mem or flush, write-thru: main mem atomic
wr-miss policies: No-write-alloc: only to main mem wr, no \$ wr.
 Write-alloc: fetch into \$, avoid wr to DRAM unless evicted
Inclusion / Exclusivity: Inclusion ($L_1 = L_2$), Exclusivity ($L_1 \neq L_2$)

MESI to reduce broadcast
 if local wr in S state (we might be the only ones) ↑ VS.



Sharing status: Dir based: block of mem (centralized - SymMP) (distributed : SGI, diff nodes)

Scoping: \$l/m, monitor bus; needs totally ordered broadcast network, can go w/o with bus

Memory access: producer-consumer: wait, use bit to indicate → OoO may give stale data

Mutual Exclusion: one process at a time, use lock

row ~~(ex)~~ coherence issue $ao = al = A$ preserve program order: R after W should be up-to-date
~~old~~ 3: $W \rightarrow 0, (a0) \rightarrow SW \rightarrow 0, (a0)$ 1, 2, 3 write serialization: 2 W + same loc. by any 2
 1: $(W, F1, (a1))$ ~~dd row~~ ps are globally seen in 0 by all processors

Eventuality: read by processor from location X that follows a wr by another processor to X returns wr if:

(1) Read & wr sufficiently separated in time (2) No other wr occur b/w R & W to X

Full bit vector scheme: ~~Bad scaling~~ each dir entry has state of cache line & bit vector w/ one share bit per 4 dir states → dir bits per cache line? ~~(ex)~~ 128 cores w/ priv \$ → 2 bit ($2^2=4$), 128 bits (bit vec) + 2

Hierarchical bit vector scheme: agr. cores into grp → map as single bit, invalids sent to all in grp

~~(ex)~~ 1024 core 64 byte cache lines, # cores/grp to reduce dir stat by 10X of phys mem?
 $2 + 1024/N \leq 64 \times 8/10 \rightarrow 2 + 1024/N \leq 512/10 \rightarrow N \geq 21 \rightarrow$ at least 21 cores

* Must store dir bits for every line in mem even if not \$d → hierarchical dir struct.

↳ or add last level cache (LLC): stores only recently used; inclusive of smaller \$

False sharing: 2+ threads modify data that happens to reside on same \$ line; interleaves access each thr. May invalidate \$ line for other thr.

~~(ex)~~ sharing \$ line w/ single bit; even if val didn't change its force to be I

Mem consistency across all rem addr; even if no cache, important for Multi-thr.

Sequential Consistency (SC): order-preserving PI: A, B p2; C, D → BA DC not allowed

* Most real time not SC → slower, needs complex L1/W, reads TSO to work

Store buffer optimization: CPU → ST buff → Shared mem; later loads can go ahead of st if diff

Total store ordering (TSO): relaxes W → R (R can be visible before W)

* Allow use of stub buff / sd → sd could take ∞ time

Weak multi-copy atomic mem: all processors see writes by another processor in some order (relaxes everything outside but internally, everything is in order) "pt of visibility"

NON-multi-copy ...: write can be partially visible to all cores (intermediate buffer)
 ALKA hierarchical shared buffering, relax all, maintain inner order

Weak: simpler to achieve high perf, easy HW, expose HW to SW | **Strong**: easy on ISA, high complexity

Fences

producer → flag data → consumer

`SW xdata(xdatahp)`
`l1 xflag, 1`
`fence r/r`
`SW xflag, (xflagp)`

Spin: l1w xflag (xflag)
`beg z xflag spin`
`fence r/r`
`lw xdata(xdatahp)`

Relaxed mem models:

$X \rightarrow Y$: X must be visible before Y is visible

$R \rightarrow R, R \rightarrow W, W \rightarrow R, W \rightarrow W$

Release consistency: relax all four

★ Compilers can reorder → bad performance

initial state	other cached	ops	Actions by this cache	Goal state	this cache	other caches	mem
Clean Exclusive ↳ CPU wr/rd = self processor issues, write hit/miss	W/S Signals	none	none	CR	yes		yes
		CPU read	none	CE	yes		yes
		CPU write	none	OE	yes		
		replace	none	I			yes
		CR	none or CCI!	CS	yes	yes	yes
		CRI	none or CCI!	I		yes	
		CI	none		impossible (if shared)		
		WR	none		impossible many...		
		CWI	none	I			yes

ID	Event	Message Shared	Cache0 State	Cache1 State	Cache2 State	Main Memory Up-to-Date?
0	CPU0: read A	0:CR	OE	I	I	Yes
1	CPU2: write B	2:CR!	I	I	OE	No
2	CPU1: read B	1:CR; 2:CCI	I	OS	OS	No
3	CPU1: write B	1:CI	I	OE	I	No
4	CPU2: write A	2:CR; 1:CCI	I	I	OE	No
5	CPU0: write B	0:CRI; 2:CCI	OE	I	I	No

Guest lectures

Apple

- Branch predictions
- Cost of mispredict grows w/ deeper pipes
- Static (simplest)
 - ↳ backward (loops) usually taken
 - ↳ fwd J is more unknown
- Dynamic (history)
 - ↳ Bimodal predictor: FSM indexed w/ PC
- Predictor aliasing:
 - 2 diff branches w/ diff PC & diff bias share same predictor index
 - ↳ Agree predictor: BTB + tags
 - ↳ Discreed: odd # faults; majority vote
 - ↳ Perception (MC)

Prefetching Occupancy = $\frac{\text{latency}}{\text{throughput}}$
 $\text{Throughput} = \frac{\text{window size}}{\text{latency}}$

- next line
- Instr stride (use PC to index into PF Table)
- Irregular (concurrency)
- Content directed (execution-based)
- Hybrid: use mul. PF to cover patterns, but complex BW intensive
- Feedback: throttle

PF using performance; change loc of b insert

WSC

- WSC: single org, homog, common layers / SW platforms, in-house services, renting VMs, request live parallelism
- Parts: Server, racks, Top of R returning switch, Accelerators (demand for ML) GPUs, TPUs, FPGAs, RCU (cyanide)
- Networking: hard to fix bw
 - ↳ N ports for N blocks v.s. fat tree, Clos @ Jupiter DC fabric → rack switch → DCAAM
- Power: transformers, UPS, PDUs
- Uninterruptible power Sys (UPS):
 - generator, remote sys / distant fresh air box
- WSC cooling: hierarchy of loops
 - Open loop: new medium replace.
 - Close loop: re-circulate, transfer efficiency = compute / total energy = $(\frac{1}{PUE}) \times (\frac{1}{SPUE}) \times (\frac{\text{Comp}}{\text{total Energy + comp}})$

PUE = facility power / IT equip power
 SPUE = server power conversion efficiency
 (c) server arch eff. (a) facility eff.

Fault tolerance:

- ↳ Power/cooling: hierarchy, scheduler
- ↳ Net.: redundancy, BW
- ↳ Scheduling: replication, dist

AWS

- Nitro: more cloud Server aka virtualize
- ↳ server CPU not shard, no limit on inst. size, no tradeoff of network/storage / custom instances
- Scale I/O workload w/ high durability & low-latency
- Like migration
- Security
- Scalability (decouple VA functions from HW), easier mgmt