



### DATA HAZARDS:

- Read - after - write dependency: The destination register for add/nor/lw is read from in add/nor/lw/sw/beg
- register reads happen in Stage 2
- register writes happen in Stage 5
- At most, next 2 instructions after add/nor/lw can trigger a data hazard in 5 stage pipeline
- RECALL:** DestReg is destination for add/nor [IMPORTANT]  
Reg B is destination for LW
- Approaches:**
  - ① **Avoid:** Pad instructions with noops ( $\leq 2$ )  
Done by user in code
  - ② **Detect + Stall:** Pad instructions w/ noops, but do it in the hardware.  
Stall IF, ID stages until front instr. hits WB  
Increases CPI - 2-3 cycles for add/nor/lw
  - ③ **Detect + Forward:** Once value is computed, forward to decode stage as needed.
    - For ADD/NOR, computed in stage 3, so can be internally-fudged (no stall needed)
    - For LW, computed in stage 4 (MEM), so can be fudged only if instr is 2-before.  
Still need stall for 1-after instruction.

### CACHE VOCAB:

- Things in cache line:  
→ tag → LRU-tracking → Valid bit → Dirty Bit
- WRITE-THROUGH:** Memory is written to whenever some data is stored
- WRITE-BACK:** Memory is written to only on eviction. This requires a dirty bit.
- Allocate-on-WRITE:** If an address used in a store is not in the cache, bring data into cache.  
→ no allocate-on-write would write directly to mem
- Fully-Associative:** Addresses can go anywhere in the cache. Same as n-way set associative cache with n cache lines (1 set).
- DIRECT MAPPED:** Every address maps to exactly one cache line. Some as 1-way set associative cache (n sets w/ n cache lines)
- Set - Associative:** Every address maps to a set of cache lines. n-way set associative → n cache lines in each set.  
→ balances latency to get data w/ hit rate
- |                       |              |                      |
|-----------------------|--------------|----------------------|
| Tag                   | SET IDX      | Block Offset         |
| leftover tag (# sets) | log (# sets) | log (# addy / block) |
- Block Size:** Amount of data in one line of cache (per tag). Increasing block size reduces overhead b/c more data is stored per line and tag also takes advantage of SPATIAL LOCALITY  
→ # blocks (# lines) = cache size / block size
- # of ways: def doesn't have to be power of 2
- block size, # sets: would be a pain if not powers of 2, so no

### BYTE - ADDRESSABLE vs. WORD - ADDRESSABLE:

64-bit ISA → 64-bit word

32-bit ISA → 32-bit word

Changes block offset size calculation; if block is  $256 = 2^8$  bytes in byte-addressable ISA, then BO = 8 bits. If it were a 64-bit word addressable ISA, then BO = 5 bits b/c 32 words in 256 bytes.

LC2K is 16-bit address, word (4 bytes) addressable.

CPI = Base + Control Hazard Stalls + Data Hazard Stalls + I-Cache Stalls + D-Cache Stalls

# LRU Bits =  $\log(\# \text{ cache lines})$

Often have 3 levels of caches = L1, L2, L3

### CONTROL HAZARDS

- Triggered by BEQ instructions: don't know what to execute next until BEQ executed (ss)
- APPROACHES:**
  - ① **Avoid:** This is dumb and impractical
  - ② **Detect + Stall:** Keep instructions in IF stage until branch is predicted
    - need instr in IF stage
    - get ALU result in EX stage, but no forwarding
    - so, stall for  $\leq 3$  instructions (IF-MEM align)
  - ③ **Speculate + Squash:** Assume taken/not-taken
    - if correct, no change needed;
    - if wrong, "squash" (replace w/ noops) all instr that are in prev pipeline stages, and load in correct PC
    - works b/c no perm change till WB/MEM
    - send noops to decode/ex/mem (IF/ID, ID/EX, EX/MEM registers)
    - 3 stalls if wrong, 0 if correct

### CACHE DESIGN CONSIDERATIONS:

- Instruction Cache (IC) vs. Data Cache (DC):**
  - can be separate or same
  - Pro: can optimize cache design for the diff use cases; e.g., IC wants higher block size because instructions have more spatial locality
  - Con: can't dynamically tradeoff between data space and instruction space in software; may be worse on some programs.
- Fully-associative caches can be slow b/c parallel tag searches are expensive.

### 3C PROBLEMS

- 3 types of cache misses exist:
  - Compulsory:** First reference to an address. These can't be completely avoided, but reduced by increasing block size (spatial locality).
  - Capacity:** Would've had a hit w/ a large enough cache; reduce by increasing cache size.
  - Conflict:** Would've had a hit w/ a fully-associative cache; reduce by increasing associativity (larger set size).
- To classify misses, simulate 3 caches:
  - ① Fully-associative cache of infinite size, but same block size as original.
  - ② Fully-associative cache w/ all other parameters same as original.
  - ③ Actual intended cache
- Miss in ①, ②, ③ → Compulsory  
Miss in ②, ③ → Capacity  
Miss in ③ → Conflict

### CACHE PARAMETERS VS. HIT RATE

- Cache Size:** Total data capacity. Bigger can take advantage of temporal locality better, but may increase cache lookup latency, degrade critical path

- Hit Rate:** Working set size / Cache Size
- Working Set:** Whole set of data executing references within a time interval
- Block Size:** Amount of data associated with a single tag. If too small, don't take advantage of spatial locality and lead to compulsory misses. If too big, transfer a lot of useless data — this both increases latency and kicks out useful data.

- ASSOCIATIVITY:** How many blocks map to the same set? Larger associativity → lower miss rate, but diminishing returns. Smaller associativity is lower cost, faster hit time

### BRANCH PREDICTION:

- Can we improve branch prediction at fetch stage?
  - branch direction (taken / not-taken)
  - branch target address
- Static Prediction:** At compile time
  - based on analysis of program
  - common: BTTFN (backwards taken, forwards not)
  - Loops are back-branches
  - Others are forward
- Overall: 60-70% of branches are taken
- Dynamic:** @ run-time
  - local predictor: for that specific line
  - global predictor: for All branches
  - 1-bit counter: predict last time
  - 2-bit: SNT ← WNT ← WT ← ST

### VIRTUAL MEMORY

- Cannot directly use virtual addresses (used by ISA) to get data;  $2^{64}$  bytes of memory is huge, and need to separate processes. Need to translate virtual addresses to physical addresses
- Use disk as "extra" space if main memory exhausted (swap partition).
- Divide memory into pages; size of virtual page is size of physical page. Map pages.
- Page Tables:** One per process, maintained by OS
  - In main mem, address already known
  - Indexed by VPN; store PPN, valid bit → NOT A TAG; space exists for all VPNs
  - Page fault: Page is on disk or invalid → Disk is evicted, invalid is unassigned
  - Page Table Entry contents:
    - PPN, Valid bit, Main Mem or Disk, Access Permissions, Dirty Page, LRU data
    - Sometimes LRU expensive, so OS uses heuristics like "accessed bit", cleared occasionally
  - Page Table Entry contents:
    - PPN, Valid bit, Main Mem or Disk, Access Permissions, Dirty Page, LRU data
    - Sometimes LRU expensive, so OS uses heuristics like "accessed bit", cleared occasionally

### HIERARCHICAL PAGE TABLES

- All entries in single-level page table are allocated when process starts; this is very long and problematic — so how do we allocate on demand?
- Have first-level page table point to second-level page table, which is created if needed — idea of a degree → generalize to n-level page table → more levels mean less memory but more lookup time



### ADDRESS TRANSLATION

- Going thru N-level page table needs N memory accesses, which is very slow
- Solution:** Translation Look-aside Buffer (TLB)
  - VPN to PPN, 16-512 entries
  - optimized cache for page lookups: low block size, only caches actual phys addrs
- Virtually-addressed cache vs. Physically addressed Cache**
  - Virtual: Faster access (no TLB / Page Table lookup in the middle)
  - can skip address translation
  - context-switching (changing processes) is expensive w/ VA cache as cache needs to be cleaned / reloaded
  - Typically: L1 is VA; L2, L3 are PA

### INTERNAL FORWARDING

- Assume register file can return updated value when read from/written to in some cycle

E.g.: ADD on line 1 to R3  
ADD on line 4 using R3 is NOT a problem b/c Line 4 is in Decode while Line 4 is in Writeback

• 21-14: REG A • 18-16: REG B

• 2-0: REG DEST • 15-0: offset field

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

### LC2K ISA : 32 bit

- 0. ADD RA RB R-Dest → R-Dest = RA + RB
- 1. NOR RA RB R-Dest → R-Dest = !RA | RB
- 2. LW RA RB #offset → RB = RA + offset
- 3. SW RA RB #offset → [RA + offset] = RB
- 4. BEQ RA RB #offset → Branch to PC + offset + 1 iff RA == RB (else +PC)
- 5. HALT RA RB →
- 6. NOOP →
- 21-14: REG A • 18-16: REG B
- 2-0: REG DEST • 15-0: offset field
- ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• 21-14: REG A • 18-16: REG B

• 2-0: REG DEST • 15-0: offset field

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

• ALL UNUSED BITS: 0 → 0-31: Least → Most Sig

•

