

## 1 CPU Power Consumption

Chip Power Consumption =  $C * Vdd^2 * F$   
 $C$  = capacitance,  $Vdd$  = voltage,  $F$  = frequency

**Dennard scaling:** Transistors get smaller but their power density stays the same. The supply voltage of a chip can be reduced by  $0.7x$  every generation

$$Power = C * Vdd^2 * F_{0 \rightarrow 1} + Vdd * I_{leakage}$$

Leakage gets worse with samller devices and lower  $Vdd$ , it also gets worse with higher temps

**Amdahl's Law:** We want to make the common case efficient, given an optimization  $x$  that accerlerates fraction  $f_x$  of the program by a factor  $S_x$ , the overall

$$speedup = \frac{1}{(1-f_x) + \frac{f_x}{S_x}}$$

## Performance

**Latency:** how long it takes to do a task

**Throughput:** total work done per unit time

$$ExeTime = \frac{Instrs}{Program} \frac{Clockcycles}{Instr} \frac{Sec}{ClockCycle}$$

**Relative Performance:** define performance as  $= 1/ExecutionTime$  "X is n times faster than Y" means  $Perormance_x/Performance_y$

## 2 Instruction Set Architecture

CISC	RISC
Emphasis on Hardware	Emphasis on Software
Multi-cycle complex instructions	Simple (single clock) instructions
Memory-to-Memory	Register-to-Register
load/store incorporated in instr.	Separate load/store instructions
Small code size	Large code size
High CPI	Low CPI
Low clock frequency	High clock frequency
Variable length instructions	Same length instructions
Complex instruction decode	Simple instruction decode
HW difficult to implement	HW easy to implement

R-type
<div>funcn7rs2rs1funct3rdopcode</div>

Opcode: basic operation of instruction (7)  
Rs1: Register source 1 operand (5)  
Rd: Register destination operand (5)  
Funct3: additional opcode field (3)  
Funct7: additional opcode field (7)

Question: Why did RISC-V only define 32 registers?

I-type
<div>immediate[11:0]rs1funct3rdopcode</div>

Question: What is an immediate?

S/B-type
<div>imm[11:5]rs2rs1funct3imm[4:0]opcode</div>

Question: Why is the immediate split?

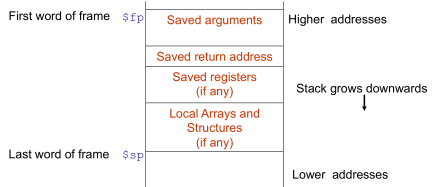
U/J-type
<div>immediate[31:12]rdopcode</div>

Constant == immediate == literal == off-set

There are 32 registers in RISC-V **ld dst, offset(base)** the base is the starting address of the array the offset is the index. When using switch statement we can use a jump table, a jump table holds addresses in memory of where the code for the jump targets are.

**Stack:** The stack is allocated in frames, it stores the state of a procedure for a limited time, the callee returns before the caller does. The things which can be saved

the stack are: local arrays, return addresses, saved registers, and nested call arguments



Callee Saved	Caller Saved
Saved registers (\$s0-\$s7)	Temporary registers (\$t0-\$t9)
Stack/frame pointer (\$sp, \$fp, \$gp)	Argument registers (\$a0-\$a3)
Return address (\$ra)	Return values (\$v0-\$v1)

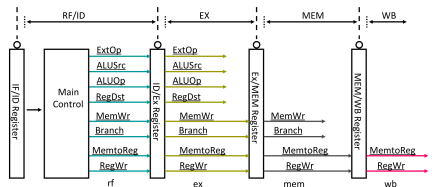
- Callee saved registers (preserved for caller)
  - Save register values on stack prior to use
  - Restore registers before return
- Caller saved registers (not preserved for caller)
  - Do what you please and expect callees to do likewise
  - Should be saved by the caller if needed after procedure call

## Procedure Call Steps

- Place parameters in a place where the procedure can access them
- Transfer control to the procedure
- Allocate the memory resources needed for the procedure
- Perform the desired task
- Place the result value in a place where the calling program can access it
- Free the memory allocated in (3)
- Return control to the point of origin

## 3 Pipelining

In pipelining we overlap instructions in different stages



There are hazards, these include **structural hazards** where a required resource is busy, **data harzards** where we must wait previous instructions to produce/consume data, and **control hazards** where next PC depends on previous instruction.

## Structural Hazards

two instructions are trying to use the same hardware within the same cycle, to solve this we can make all the instructions the same length

## Data Dependencies

Dependencies for instruction  $j$  following instruction  $i$

- Read after Write (RAW or true dependence)

Instruction  $j$  tries to read before instruction  $i$  tries to write it

- Write after Write (WAW or output dependence)

Instruction  $j$  tries to write an operand before  $i$  writes its value

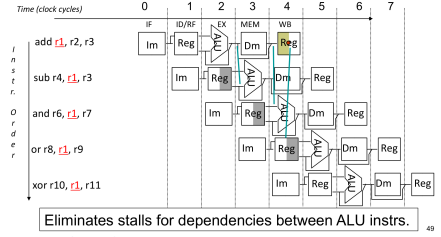
- Write after Read (WAR or (anti dependence))

Instruction  $j$  tries to write a destination before it is read by  $i$

**Solutions for RAW Hazards:** We can delay the reading of an instruction until data is available, to do this we can insert pipeline bubbles, can also write to the register file in the first half of a cycle and then read in the second half.

**Forwarding:** Another solution is forwarding or pushing the data to an appropriate unit.

We can also reorder instructions to deal with RAW hazards



## MEM hazard

- if (MEM/WB.RegWrite and (MEM/WB.RegisterRd  $\neq$  0) and not (EX/MEM.RegWrite and (EX/MEM.RegisterRd  $\neq$  0) and (EX/MEM.RegisterRd == ID/EX.RegisterRs)) and (MEM/WB.RegisterRd == ID/EX.RegisterRs))  
ForwardA = 01
- if (MEM/WB.RegWrite and (MEM/WB.RegisterRd  $\neq$  0) and not (EX/MEM.RegWrite and (EX/MEM.RegisterRd  $\neq$  0) and (EX/MEM.RegisterRd == ID/EX.RegisterRt)) and (MEM/WB.RegisterRd == ID/EX.RegisterRt)) and (MEM/WB.RegisterRd == ID/EX.RegisterRt))  
ForwardB = 01

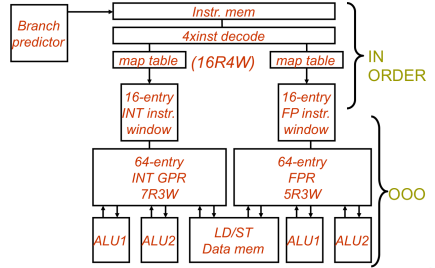
## Control Hazards

A control hazards is like a data hazard on the PC, we cannot fetch the next instruction if we don't know the PC

Some solutions for control hazards are stalling on branches, predicting taken or not taken. We need to flush the pipeline if we predict wrong, in a 5-sage pipeline we only need to flush 1 instruction

## 4 Out of Order Execution and ILP

We want to avoid in-order stalls so we use out of order execution to re-order instructions based on dependencies



A **superscalar processor** is a CPU that implements a form of parallelism called instruction-level parallelism within a single processor. I.e we can launch multiple instructions every cycle.

There are some issues with multiple instructions executing at once, we need to double the amount of hardware, we introduce hazards, branch delay, & load delay. We can rename (map) architectural registers to physical registers in decode stage to get rid of false dependencies

Superscalar + Dynamic scheduling + register renaming

```
add $t0, $t1, $t2    sub $t0, $t1, $t2
or $t3, $t0, $t2     and $t5, $t0, $t2
```

There are some limits to ILP and pipelining:

- Limited ILP in real programs
- Pipeline overhead
- Branch and load delays exacerbated
- Clock cycle timing limits
- Limited branch prediction accuracy (85%-98%)
- Even a few percent really hurts with long/wide pipes
- Memory inefficiency
- Load delays + # of loads/cycle

## 5 Caches

	Access time	Capacity	Managed by
Registers	1cycle	~500B	software/compiler
CPU Chip			
Level 1 Cache	1-3cycles	~64KB	hardware
Level 2/3 Cache	10-20cycles	1-100MB	hardware
Chips			
DRAM	~100cycles	~100GB	software/OS
Mechanical Devices			
Disk	10 <sup>6</sup> -10 <sup>7</sup> cycles	~2TB	software/OS
Tape			

**Principle of locality** Programs work on a relatively small portion of data at any time, so we can predict data accessed in near future by looking at recent accesses. There are two types of locality: **spatial** and **temporal**, **spatial** locality is if an item has been accessed recently, nearby items will tend to be referenced soon, and **temporal** is if an item has been referenced recently, it will probably be accessed again soon.

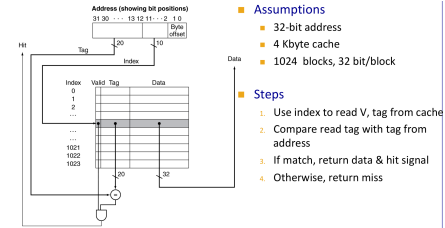
How is data stored in a cache? There are three ways:

- Direct mapped (single location)
- Fully associative (anywhere)
- Set associative (anywhere in a set)

## Direct Mapped Cache

Location in cache determined by (main) memory address, we use the lowest order bits to determine this address.

To determine which particular block is stored in a cache location, we look at **tag** and **valid** bits, the **tag** bits are the high-order bits of the memory address, the **valid** bit represents **1** = present, **0** = not present



- Block** – Minimum unit of data that is present at any level of the hierarchy
- Hit** – Data found in the cache
- Hit rate** – Percent of accesses that hit
- Hit time** – Time to access on a hit
- Miss** – Data not found in the cache
- Miss rate** – Percent of misses (1 - Hit rate)
- Miss penalty** – Overhead in getting data from a higher numbered level
  - Miss penalty = higher level access time + Time to deliver to lower level + Cache replacement / forward to processor time
  - Miss penalty is usually much larger than the hit time
  - This is in addition to the hit time
- These apply to each level of a multi-level cache**
  - e.g., we may miss in the L1 cache and then hit in the L2

## Average Memory Access Times:

$$AccessTime = hittime + missrate \times misspenalty$$

## 3 C's of Cache Misses:

- Compulsory** – this is the first time you referenced this item
- Capacity** – not enough room in the cache to hold items  
this miss would disappear if the cache were big enough
- Conflict** – item was replaced because of a conflict in its set  
this miss would disappear with more associativity

$$CPIpenalty = missrate \times misspenalty$$

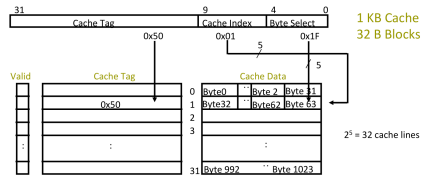
When we encounter a miss the pipeline stalls for an instruction or data miss

## Cache Blocks

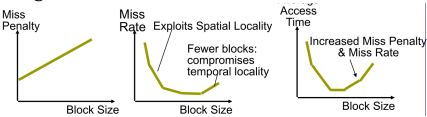
Assuming a  $2^m$  byte direct mapped cache with  $2^m$  byte blocks

- Byte select – The lower m bits
- Cache index - The lower (n-m) bits of the memory address
- Cache tag - The upper (32-n) bits of the memory address

Direct Mapped Problems: **Thrashing** If accesses alternate, one block will replace the other before reuse



**Block Sizes:** Larger block sizes take advantage of spatial locality, however it also incurs larger miss penalty since it takes longer to transfer the block into the cache.

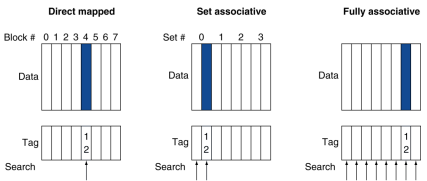


### Fully Associative Cache

Opposite extreme in that it has no cache index to hash, it uses any available entry to store memory elements, there are no conflict misses, only capacity misses, and we must compare cache tags of all entries to find the desired one

### N-way Set Associative Cache

Compromise between direct-mapped and fully associative, each memory block can go to one of N entries in cache, and each set can store N blocks; a cache contains some number of sets



### Tag & Index with Set-Associative Caches

Given a  $2^n$  byte cache with  $2^m$  byte blocks that is  $2^a$  set-associative, the cache contains  $2^{n-m}$  blocks, and each cache way contains  $2^{n-m-a}$  blocks, and the cache index is  $n - m - a$  bits after the byte select

Organization	# of sets	# blocks / set	12 bit Address
Direct mapped	16	1	tag index blk off 4 4 4
2-way set associative	8	2	tag index blk off 5 3 4
4-way set associative	4	4	tag ind blk off 6 2 4
8-way set associative	2	8	tag i blk off 7 1 4
16-way (fully) set associative	1	16	tag blk off 8 4

Some cons with associative caches are that there is an area overhead and more latency

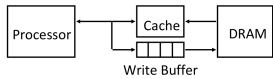
The **replacement policy** for a N-way set

associative caches choosing the least recently used thing

### Cache Write Policies

- Write-through (write data go to cache and memory)
  - Main memory is updated on each cache write
  - Replacing a cache entry is simple (just overwrite new block)
  - Memory write causes significant delay if pipeline must stall

- Write-back (write data only goes to the cache)
  - Only the cache entry is updated on each cache write so main memory and the cache data are **inconsistent**
  - Add “dirty” bit to the cache entry to indicate whether the data in the cache entry must be committed to memory
  - Replacing a cache entry requires writing the data back to memory before replacing the entry if it is “dirty”



#### Use Write Buffer between cache and memory

- Processor writes data into the cache and the write buffer
- Memory controller slowly “drains” buffer to memory

#### Write Buffer: a first-in-first-out buffer (FIFO)

- Typically holds a small number of writes
- Can absorb small bursts as long as the long-term rate of writing to the buffer does not exceed the maximum rate of writing to DRAM

### Write Miss Options:

Steps	Write through				Write back	
	Write allocate		No write allocate		Write allocate	
	fetch on miss	no fetch on miss	write around	write invalidate	fetch on miss	no fetch on miss
1	pick replacement	pick replacement		invalidate	pick replacement	pick replacement
2				invalidate tag	[write back]	[write back]
3	fetch block				fetch block	
4	write cache	write partial cache			write cache	write partial cache
5	write memory	write memory	write memory	write memory		

**Splitting Caches** Most chips have separate caches for instructions and data, some advantages are that we have extra bandwidth and low hit time, but miss rate will be higher.