

grai adds sign bit into empty bits

Two complement Tricks

$$x + \sim x + 1 = 0$$

$$\text{MAX_INT} + 1 = \text{MIN_INT}$$

$$-1 = 0x\text{FF...F}$$
 (all 1's in binary)

$$\text{MAX_INT} = 2^{n-1} - 1, \text{MIN_INT} = -2^{n-1}$$

$\sim x + 1$ to negate

Prologue:

```
addi sp, sp, -12    lw ra, 8(sp)
sw $t0, 0(sp)      lw $t1, 4(sp)
sw $t1, 4(sp)      lw $t0, 0(sp)
sw ra, 8(sp)       addi sp, sp, 12
```

Epi logue:

```
struct stack * stk = (struct stack *) malloc(sizeof(struct stack));
```

```
if (!stk) {
```

```
    return NULL;
```

```
    stk->size = size;
```

```
    stk->topIndex = -1;
```

```
    stk->stackArray = (int *) malloc(sizeof(int) * size);
```

```
    if (!(stk->stackArray)) {
```

```
        free(stk);
```

```
        return NULL;
```

```
}
```

Branch offset (in bytes)

Jump offset (in bytes)

```
return stk;
```

```
}
```

mod and read from right to left

$$\begin{array}{lll} 2^5 = 32 & 2^6 = 64 & 2^7 = 128 \\ 2^8 = 256 & 2^9 = 512 & 2^{10} = 1024 \end{array}$$

typedef → don't have to specify struct
and can't be overwritten

array var is a pointer to the 0th elem

null terminate strings when copying/mallocing

have to keep org malloc ptr to use free()

cannot return pointer to things on the stack

```
int push(int x, struct stack * stk) {
```

```
    if (stk->topIndex == stk->size - 1) {
```

```
        int *stackArray = (int *) realloc(stk->stackArray,  
                                         stk->size * 2);
```

```
        if (!stackArray) return 0; }
```

```
        stk->stackArray = stackArray;
```

```
        stk->size = stk->size * 2;
```

```
}
```

```
stk->topIndex++;
```

```
stk->stackArray[stk->topIndex] = x; }
```

a0, a1 → return value

a0-a7 → args

vals to save b4 jal → a0-a7, t0-t6, ra

vals to restore b4 jalr → sp, s0-s11

strings in func should be static or
malloced if need to be returned

Another option is to allocate on the stack in
the caller

Callee: ensure saved registers and sp
are the same

Caller: must save other registers if needed

Stack: local var, const (defined in func)

Heap: malloc

Static: static var, global var, string literals, const^{not in} func

Code: const (DEFINE), machine instructions

R-type add rd rs1 rs2

I-type addi rd rs1 imm J-type (load) lw rd, imm(rs1)

S-type sw rs2 imm(rs1)

SB-type beq rs1, rs2, label

imm depends on offset (= mul of 4 bytes)

Loop : beg x19, x10, End

addi

addi

j Loop

End:

4 inst x 4 bytes = 16 = imm

1
2
3
4

Packet* read_into_packet(size_t max_size) {

int size = 0;

char buf[4];

Packet* pkt = malloc(sizeof(Packet));

pkt->data = malloc(sizeof(char) * (max_size));

while (size < max_size) {

int n = read4(buf);

if (!n) break;

int i = 0;

while (i < n) {

pkt->data[size + i] = buf[i++];

y

pkt->size = size;

return pkt;

(can reach
same with
8 inst
of PC)

UJ-type jal x0, label
or j label

Jalr ret = jr ra = jalr x0, ra, 0

To call func at any 32 bit abs address:

lui xl, <hi 20 bits>

jalr ra, xl, <lo 12 bits>

Procedure

Associativity

1 ++ -- suffix left-to-R

2 ++ -- prefix Right-to-L

*

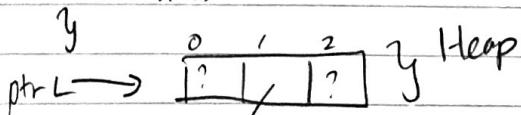
int main()

char **ptr = malloc(4u);

ptr[0] = "c";

printf ("%c", ptr[0]);

return 0;



o l
|'c'l'o'| ← str literals
in static

output: c

Jump Branch (offset in bytes)

```

for (i = 0; i < len(A); i += block) {
    for (j = 0; j < len(B); j += 1) {
        for (k = 0; k < block; k++) {
            if (k + i) < len(A) {
                fn(A[i * block + k], b[j])
            }
        }
    }
}

```

Branch Predictor: N-entry, direct-mapped

- when inst is branch check br! for $pc[7:2]$ is set, set next pc to $pc + \text{offset}$ (predicted-taken)

when err!, if br taken but pred not, kill

if br not taken but pred taken, lcn!

Entry starts at 01, if taken add 1, if not sub 1,

if upper bit is 1, assume taken

return target location

small stack to check for jalr and jal

when writing to ra, put ret4 to stack

when reading from ra, pop stack

multiple processor reading Ok

Processor A wants to write, write-allocate,

set the dirty bit, and broadcast, others

processors invalidate the entry if they have it.

Miss in upper level, if no one has it, just fetch

If a processor has it w/ dirty bit set, flushes

entry (write to mem), clears the dirty bit and

gives the vcu to the requesting processor

coherence misses, caused by writes

by two processors on the same data

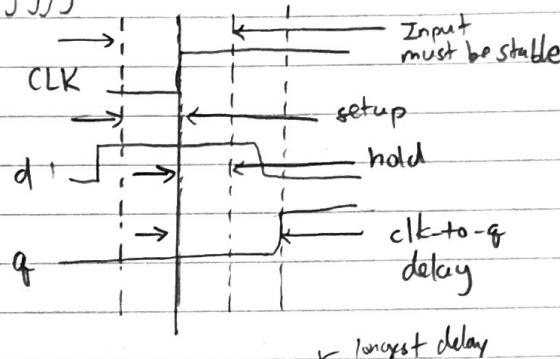
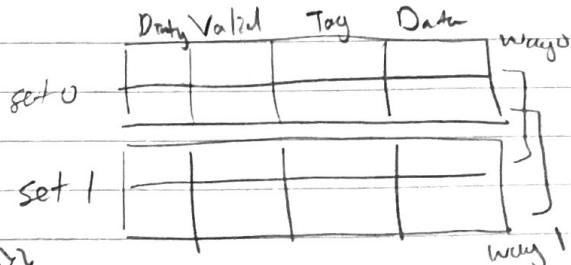
Incoherence misses - two processes have

diff working sets \rightarrow reduced cache

size and miss rate

To solve = use virtual mem, great if

we never have to fetch from disk.



$t_{clk} > t_{clk-to-q} + t_{logic} + t_{setup}$

$t_{clk-to-q} + t_{logic}/shortest > t_{hold}$

Latency = Time for 1 inst to finish

SS to Pipelined, latency \uparrow

SS = add up delay for all stages

Pipelined: # stages \times clock period

Throughput: # of inst processed/second

SS to Pipeline = throughput \uparrow

stages / latency

Branch - kill next 2 inst

load - stall 1 cycle

data w/o forward - 2 stalls

Cache Write Hit

Write back, or Write thru

Write Miss

Write allocate - fetch then write

No wa - write to mem and don't fetch

1. How many accesses / loop
 - a. count write and read
 - $a[i] += 5$ is read & write
2. How many loops/steps per block
3. How long is cache compared to array
 - a. what % of array can fit
 - b. Can we reuse cache in later func
4. Do we reuse blocks once we "leave" them
5. How to fix various misses

$$AMA1 = L1 HT + L1 MR (L2 HT + L2 MR * L2 MP)$$

Global MR - fraction of all accesses that missed at that cache level
 $\frac{\# \text{ cache misses}}{\# \text{ total accesses}}$

Local MR = fraction of accesses that occur at that cache level that miss
 Ex. $L2 \text{ misses} / L2 \text{ accesses}$

CL - output as a func of input

Sequential logic

max delay = crit path

Compiler

In: foo.c Out: foo.s (Assembly lang code)

may contain pseudo-instr (assembler understands
but not machine)

Assembler: reads and uses directives

replace pseudo inst, produce Machine lang

creates object file In: foo.s Out: foo.o

Assembler directives: give directions to assembler

but do not produce machine inst.

.text (machine code).data (binary rep of data in src file) to change after the edge

.global sym (declares sym global, can be ref'd from other files) t_{hold} <= t_{clk-to-Q} + t_{shortest CL from any reg to a reg input} (min delay)

.string str (store the string str in mem & null terminate)

.word w1..wn: store the n 32-bit quantities in successive memory words t_{clk-to-Q} + t_{longest CL from reg to any input + t_{setup}}

Forward ref: jump within a file is okay (PC relative)

branch requires two parses

Symbol table: "items" in this file that others may use (.data)

Relocation table: list of "items" this file needs the address of later

Ex. external jumps or data in static section (la inst)

Object file format: header, text seg, data seg, reloc info, sym table, debugging info.

Linker: In: foo.o, libc.o Out: a.out

Combines object files into executable, enable separate

compilation of files. Takes text seg from each, then

data seg from each, then resolve abs address

PC Relative (breq, bne, jal) never relocate

External func Ref (usually jal) always relocate

Load/Store to var in static area, relative to go (relocate)

Linker knows len of txt seg and data seg and ordering, calc the

abs address of jumps and data referenced

Output: executable w/ text and data (plus header)

Loader: In: executable Out: (program is run)

load to mem from disk and start, usually the OS

flip-flop - one bit of state, sampled on rising edge

Key - several bits of state, sampled on rising edge

Boolean

$$1+A=1 \quad A+\bar{A}=1 \quad A+AB=A$$

$$0B=0 \quad B\bar{B}=0 \quad A+\bar{A}B=A+B$$

$$(A+B)(A+C)=A+BC$$

$$\overline{AB}=\overline{A}+\overline{B} \quad \overline{A+B}=\overline{A}\overline{B} \quad \text{DeMorgan's law}$$

Setup Time: input must be stable b4 edge

Hold Time: input must be stable after edge

CLK-to-Q: how long it takes the output

to change after the edge

t_{hold} <= t_{clk-to-Q} + t_{shortest CL from any reg to a reg input} (min delay)

.word w1..wn: store the n 32-bit quantities in successive memory words t_{clk-to-Q} + t_{longest CL from reg to any input + t_{setup}}

Forward ref: jump within a file is okay (PC relative)

<= CLK-period (max delay)

Improving Cache

Reduce time to hit \rightarrow smaller cache

Reduce miss rate \rightarrow bigger cache

Reduce miss penalty \rightarrow cache levels

↑ Associativity

Hit time \uparrow

↑ Entries

hit time \uparrow

Miss rate \downarrow

Miss rate \downarrow

Miss penalty mostly unchanged Miss penalty unchanged

increase in hit time might be worse

even if hit rate goes up.

↑ block size

Hit time unchanged

Miss penalty

Miss rate \downarrow then \uparrow

\uparrow with \uparrow blk size

more spatial

more conflicts

bc less blocks

but amortized over whole block

w/ fixed const init latency

Global miss rate: misses in

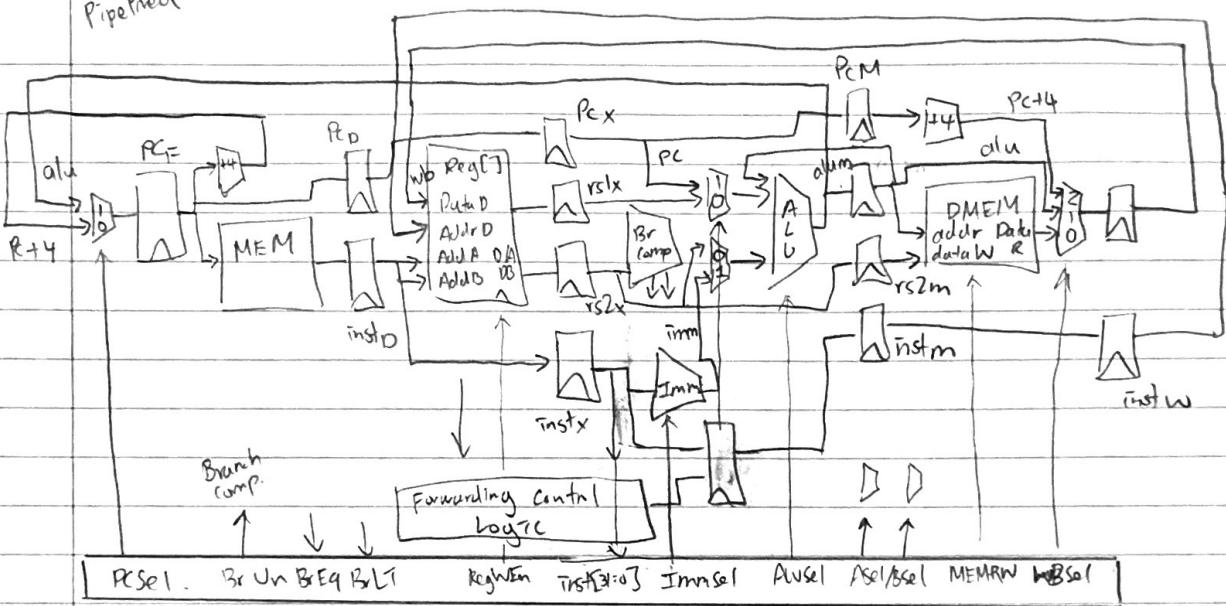
this cache / total # of mem access by CPU

local miss rate: LMR_{L2} = L2 misses / L1 misses

= L2 misses / total L2 accesses

L2 local MR \gg GMR

Pipelined



lw takes the longest Processor Performance $\frac{\text{Time}}{\text{Program}} = \frac{\text{Inst}}{\text{Program}} * \frac{\text{Cycles}}{\text{Inst}} * \frac{\text{Time}}{\text{cycle}}$
 $\text{Time/Cycle} = 1/\text{Frequency}$

Structural hazards: 2 or more inst compete for access to same resources; stall or add hardware ex. IMEM and DMEM

$$\frac{\text{Energy}}{\text{Program}} = \frac{\text{Inst}}{\text{Program}} * \frac{\text{Energy}}{\text{Inst}} = \frac{\text{Inst}}{\text{Program}} * C^2$$

$$\text{Energy from law} = \text{Power} * \text{energy efficiency} \\ (\text{task/sec}) = (\text{joules/sec}) (\text{Tasks/Joule})$$

Reg file - sep RW ports

Data hazard: value not updated before use; stall or forward

Control hazards: branch & jump

Load data hazard: must stall one inst (load delay slot)

leads to wrong inst fetched

Solution: convert fetched inst to NOPs for next two cycles
 branch prediction to reduce performance loss, or unrelated inst

Cache

$$\text{Tag} = \text{Addr size} - \text{index size} - \log_2(\# \text{bytes}/\text{block}) \quad \text{Index} = \log_2(\text{cache size}/\text{block size} * N) \quad \geq$$

$$0 = \log_2(\text{block size})$$

8 blocks

$$\log_2(\# \text{of sets})$$

Fully Associative: block can go anywhere (no index)



4 way set

Direct Mapped: block goes to one place

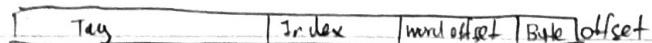
N-way set associative: N places for a block

$$\text{AMAT} = \text{Time for a hit} +$$

$$\text{miss rate} * \text{Miss Penalty}$$

Write-back: write only to cache and write to mem when evicted. (need a dirty bit)

Fixed-size cache and fixed block size

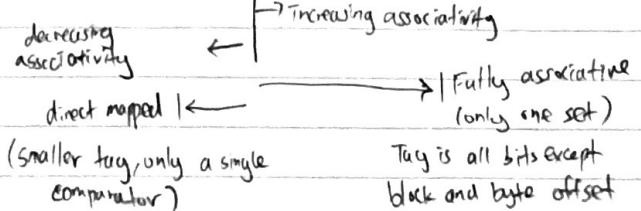


Cache hits

No write allocate: only write to mem

Write allocate: fetch into cache

Compulsory, conflict, capacity



(smaller tag, only a single comparator)

TAG is all bits except block and byte offset

Prologue Epilogue Stack: local var, const (defined in func)

```

addi sp sp -12    lw r0, $0(sp)
sw $0 0(sp)      lw $1, 4(sp)
sw $1 4(sp)      lw $0, 0(sp)
sw $2 8(sp)      addi sp sp 12

```

I-leap: malloc
Static: static var, global var, string literals, const (not in func)
Code: const (DEFINE), machine instructions

Loop beg x19 x10, end),
addi
addi
j loop
end:
4 inst × 4 bytes = 16 = imm

Typedef → don't have to specify struct and can't be overwritten

Array var is a pointer to the 0th elem

null terminate string when copying/mallocing

have to keep malloc ptr to use free

cannot return pointer to things on the stack

Strings in func should be static or malloced

If need to be returned

Another soln is to allocate on the stack in the caller

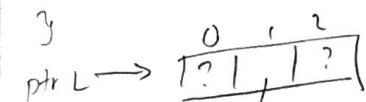
Callee: ensure saved registers and sp are the same

Caller: must save other registers if needed

```

int main() {
    char *ptr = malloc(40);
    ptr[0] = 'c';
    printf("%c", ptr[0]);
    return 0;
}

```



Output: C

strlen doesn't include the null terminator

strcpy includes null terminator

Programmed I/O

Not ideal → CPU has to execute transfer, could be doing other work, device speeds don't align well w/ CPU speeds, energy cost of using CPU when it can be replaced w/ simpler hardware

DMA
Allows I/O devices to directly R/W mem mem
DMA Engine - contains registers written by CPU

MSC
 $I=0$ for fully associative

Page table base register = # of bits of physical mem

a0, a1 → return value

a0-a7 → args

vals to save b4 jals → a0-a7, to-tb, ra

vals to restore b4 jals → sp, s0-s11

R-type add rd rs1 rs2 I-type lw rd imm(rs1) (load)

I-type addi rd rs1 imm U-type lui x10, 0x87654

S-type sw rs2 im(rs1) used 11 for large imm
SB-type beq rs1, rs2, label or -1 for the imm offset then
imm depends on offset (mul of 4 bytes)

UJ-type jal x0, label

or j label

can reach from w/l n 2^{18} inst from PC

jalr ret = jr ra = jalr x0, ra, 0

To call func at any 32 bits abs address:

lui x1, <hi 20 bits>

jalr ra, x1, <lo 12 bits>

More floating

largest finite pos value
 $(2 - 2^{-23}) \times 2^{127}$

smallest normalized 2^{-126}

step size = 2^{-23+n}

$n = x - 127$

Polling → forces hw to wait for ready bit, low latency, good when device always busy or can't move on until device replies, can't do anything while polling

Interrupts → hw raise exception when ready. CPU changes PC to execute code in the interrupt handler. Can do work while waiting, can wait on more than 1 device, good for devices that take long to respond. Non-deterministic when interrupt occurs, some overhead

Map Reduce

map → returns newly distributed datasets

flatMap → similar to map but each input

can be mapped to 0 or more outputs

reduceByKey → returns (k, v) pairs where

each key are aggregated using some reduce func

reduce → Agg regardless of key using a func

load reserved

lw rd, rs1

load the word pointed to by rs

into rd and add reservation

store condition sc rd, rs1, rs2

store value in rs2 into mem

lbranch putting tag rs1 if

the reservation is still valid and

wire status in rd

0 if success

non zero → failure, updated

by others, write doesn't

happen

Exponent	Significant	Meaning
0	Anything	Denorm
1-254	Anything	Normal
255	0	Infinity
255	nonzero	NaN

smallest rep pos # 2^{-149} (i.e., $2^{-26} \times 2^{-23}$)
 second smallest 2^{-148} (i.e., $2^{-26} \times 2^{-22}$)

Modified - up-to-date, dirty, no

other cache has a copy, can write,
 mem out of date

Owned - up-to-date, other caches may have copy
 but must be in shared state)

Shared - up-to-date, other caches may have a copy

Exclusive: up-to-date, no other cache has copy,
 ok to write, memory up-to date

Invalid = not in cache

Polling

$$\frac{\text{clocks}}{\text{sec}} = \frac{\text{polls}}{\text{sec}} * \frac{\text{clocks}}{\text{poll}}$$

% processor for polling = $\frac{\text{cpu}}{\text{clocks/sec}}$

Interrupt - caused by an event external
 to curr program running (key press)

Exception - caused by event during execution
 of inst of curr running program (div by 0)

Trap - action of servicing interrupt or exception

by hardware jump to "interrupt or trap handler" code

Interrupt is efficient for low rate

Availability = $\text{MTTF}/(\text{MTTF} + \text{MTTR}) \text{ in \%}$

Annualized Fail Rate ($1000 \text{ disks} \times 0.760 \text{ hrs/year}$)
 $\frac{1}{100,000 \text{ hrs}} = 0.76 \text{ /year}$

RAID 0

↳ Spread data across disks, bandwidth ↑ linearly.
 doesn't help w/ latency

RAID 1

↳ Mirror the data, writes limited by single disk speed,
 expensive, reads may be optimized

RAID 3

↳ Striping w/ parity disk, subtract P from sum to fix err

RAID 4

↳ Striping w/ parity disk on the block level, good for
 random reads, bad random writes (need to write to each
 P disk), good for small reads

RAID 5

↳ Striping w/ interleaved parity, independent writes possible,
 ↳ Write Algo: Read old data, parity → write new data & parity

↳ High I/O rates

Latency - Time to do one task

Bandwidth - Tasks finished per unit time

Time = $\frac{\text{seconds}}{\text{program}}$

$$= \frac{\text{Inst}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Inst}} \times \frac{\text{seconds}}{\text{clock cycle}}$$

↑
 total # of
 inst executed
 to run program

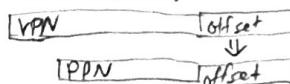
↑
 Avg # of
 cycles/inst

↑
 sys clock period
 (1/F)

PUE = $\frac{\text{total building power}}{\text{IT equipment power}}$

Energy Efficiency = $\frac{\text{Utilization}}{\text{Power}}$

Virtual Memory



PTE - 1 Valid Bit
 2 R/W Bits + PPN
 1 Dirty Bit DPN

Page Table Base Register
 ↳ points at start of PT

Disk Access Time

$$= \text{Seek time} + \text{Rotation Time} + \text{Transfer Time} + \text{Controller}$$

seek time = # of tracks / 3 * time to move
 across one track

rotation Time = $\frac{1}{2}$ time of a rev

Transfer Time = size / transfer rate

Flush has no seek time

Bit position	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Encoded data	p1	p2	d1	p4	d2	d3	d4	p8	d5	d6	d7	d8	d9	d10	d11
Coverage	p1	X		X		X		X	X		X	X	X		X
	p2		X	X			X	X			X	X		X	X
	p4				X	X	X	X				X	X	X	X
	p8								X	X	X	X	X	X	X

Two's Complement Tricks

$$X + \sim X + 1 = 0$$

$$\text{MAX_INT} + 1 = \text{MIN_INT}$$

$$-1 = 0xFF \dots F$$

$$\text{MAX_INT} = 2^{n-1} - 1, \text{MIN_INT} = -2^{n-1}$$

$$\sim X + 1 \rightarrow \text{negate}$$

$$2^5 = 32 \quad 2^6 = 64 \quad 2^7 = 128$$

$$2^8 = 256 \quad 2^9 = 512 \quad 2^{10} = 1024$$

Amdahl's Law

$$\text{Speedup} = \frac{1}{(1-P) + \frac{F}{S}}$$

↑ speed-up part
 ↓ accelerated factor

OMP - shared memory

pragmas: omp parallel for

↳ divide index per thread sequentially
 (no false sharing)

reduction - specifies 1 or more

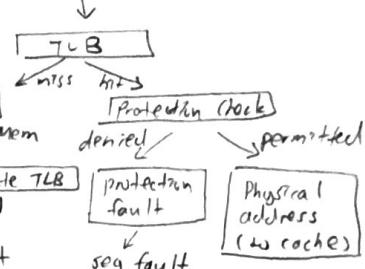
private var that combines at the end

func: (+, *, -, min, max, &, &&,
 |, ||, ^)

critical ⇒ locks, only 1 thread
 runs at a time

Each hardware thread has its own cache

Virtual Address



ECC

6 bits can cover $= 2^n - 1 - n$

Incorrect bit =

. sum(parity bit label)

Eg. P2 and P8 in error

↳ position 10 is wrong

Branch offset (in bytes)

Jump offset (in bytes)

Divide by base and

continue w/ remainder

until 0, read remainders

from last to first