

NUMBER REPRESENTATION

TYPE	RANGE	LOWEST	HIGHEST	# OF REPRESENTABLE NUMBERS
unsigned	$[0, 2^{n-1}]$	00...00	11...11	2^n
sign & magnitude	$[-(2^{n-1}), 2^{n-1}]$	11...11	01...11	$2^n - 1 \leftarrow 2 \text{ nps of } \emptyset$
2's complement	$[-2^{n-1}, 2^{n-1}-1]$	10...00	01...11	2^n
Biased	bias, $2^n - 1 - \text{bias}$	00..00	11...11	2^n

TWO'S COMPLEMENT: first bit=1 → negative

to find magnitude: ① invert bits of negative # ② add 1 ③ Negate unsigned result

shift left logical in 2's comp: multiply by 2
shift right in 2's comp: floor divide by 2

Bias: $x + b = v$
 \uparrow bias \uparrow comp-strings ($\gg p, \gg q$):
 Unsigned # you want to represent char* str1 = [*(char*)p]
 binary to represent char* str2 = [*(char*)q]

C (pass by value)

types: char: 8 bits = 1 byte int8_t = 8 bit signed
 int: 32 bits = 4 bytes

False values: '\0', 0, NULL

STRUCTS: USE `(*struct).field`
`STRUCT → FIELD`

POINTERS: USE `&` to get address
 USE `*>` to get value

pointer arithmetic: increment n by
 (ex: int = +4) `sizeof(pointer-type)`

arrays: int arr [] = {1, 2, 3};
`char* hello = "hello"; a[i] → *(a+i)`

`char* hello = malloc(sizeof(char)*6);`

strlen does NOT include NULL acct for null terminator

* Strings end w/ null terminator '\0'

MEMORY ALLOCATION:

`malloc(size_t size)`: returns void * pointer
 ↗ check if return NULL initializes w/ garbage
`free(void *ptr)` : must free any malloc'd pointer once

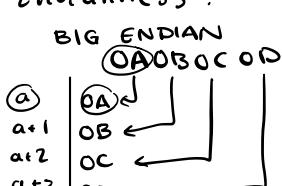
`realloc(void *ptr, size_t size)`: reallocate memory

`calloc(num_items, size)`: allocates & initializes to zero

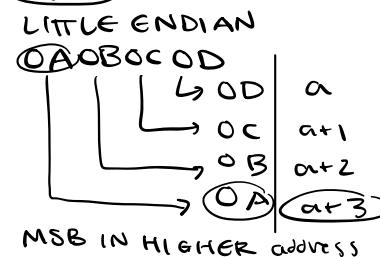
`void *memcpy(void *dest, const void *src, size_t n)`

`void *memmove(void *dest, const void *src, size_t n)`

endianness :



(RISC-V) OPOCOBAD



logical left shift: shift left add zeros = multiply pow2

logical right shift: shift right zero extend = divide unsigned by power of 2

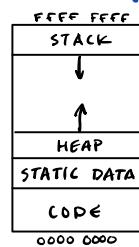
arithmetic right shift: shift right sign extend ≠ divide signed by power of 2

SLIIx12 X10 0X10 shifts LEFT by 16 bits 0X000034FF → 0X34FF0000

SRIIx12 X12 0X08 shifts RIGHT by 8 bits 0X34FF0000 → 0X0034 FF00

SRAI a0 a0 0X02 sign extend shifts RIGHT by 4 bits 0XF0000000 → 0XF000000

MEMORY MANAGEMENT



grows down (LIFO), function / local vars, erased when fn return
`int *foo = malloc(sizeof(int))` #DEFINE X
`foo = stack` ⇒ foo = heap X in code/text
 dynamically allocated (malloc, calloc, realloc), bottom up, persists over calls
 global vars outside of funcs, immutable, string literal, no grow/shrink
 loaded @ start, executable read-only code, functions, machine instructions

Floating Point

32 bit Single precision

BiAs = -127

Normalized: $\text{Value} = (-1)^{\text{sign}} \times 2^{\text{Exp} + \text{BiAs}} \times 1.\text{Significand}_2$

Denormalized: $\text{Value} = (-1)^{\text{sign}} \times 2^{\text{Exp} + \text{BiAs} + 1} \times 1.\text{Significand}_2$

EXPONENT	SIGNIFICAND	MEANING
0	Anything	DENORM
1-254	Anything	NORMAL
255	0	INFINITY
255	Nonzero	NAN

VALUE → BINARY

-47.34375 ① 47 = 10111

② 0.34375 = 0.25 + 0.0625 + 0.03125

= $2^{-3} + 2^{-4} + 2^{-5} = 0.01011$ floating point closest to 0 = denorms

③ 10111, 01011 = 1.011101011×2^5

④ $X - 127 = 5$

$X = 5 + 127 = 10000100$

1 10000100 011101011, pad w/ zeros,

RISC V

callee: s0-s11, sp

caller: a0-a7

60-t6

ra

Save & restore

PC = program counter
 → where we are in running code

JUMP instructions MODIFY PC

funct7, funct3, opcode = SPECIFIC operation

RType = arithmetic (2 source registers)

IType = immediate arithmetic (1 source, 1 immediate)

I+T-type = immediate shift (slli, srli, srai)

S-type = store

B-type = branches

U-type = upper immediate

J-type = jump (jal, j)

add memp RISC commands

lb: load byte will go to LSB & MSB smear to fill

* lb x12 i(x) = 0x93 → 0b1001...

→ 0xFFFFFFFF93

lbu: unsigned smear w/ 0

→ 0X00000093

andi with 000000FF to isolate LSB

iui: writes 8 bits to LSB and fills other 24 bits w/ 0

auipc: build 32 bit with a iui

⇒ add to this

j label
 (pseudo for jal x0 label)

jal ra label
 ↑ destination register w/ address to return to

jr ra
 (pseudo for jal x0 ra 0)

jalr ra sl 0
 (jump to 0+address in sl)

jump there, don't come back

jal ra label

jump there, come back

jalr ra sl 0

(jump to 0+address in sl)

gcc = C compiler

C → RISC V
 outputs labels, pseudo instructions, relative addressing, error checking, reorders code

RISC V → object files
 replaces pseudo instructions w/ assembly

relocation table (external labels)

symbol table (labels, static data, relative addresses)

arithmetic, logical shifts

object files → single binary executable
 put tgt data & text from .o files → put on end of text

resolve refs w/ reloc table → fill in addresses

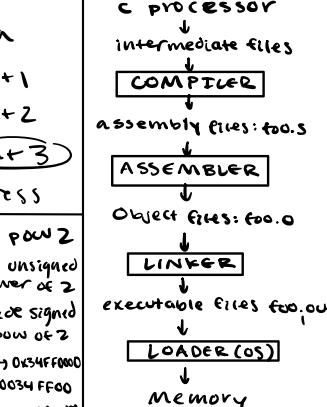
load program into memory for execution

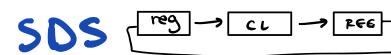
initialize registers

run a.out

part of operating system

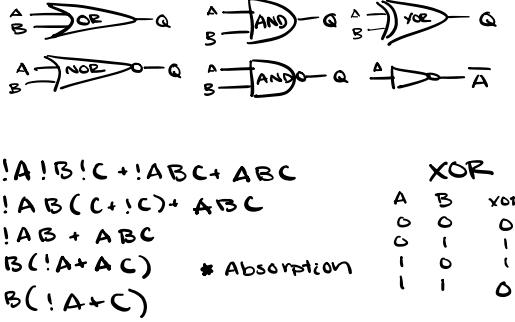
CALL



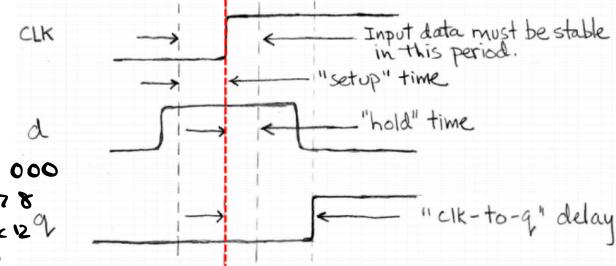


commutative	AND AB=BA	OR A+B=B+A
associative	AB(C)=A(BC)	A+(B+C)=(A+B)+C
identity	IA=A	0+A=A
null	OA=0	I+A=I
absorption	A(A+B)=A	A+AB=A
distributive	(A+B)(A+C)=A+BC	AC(B+C)=AB+AC
deMorgan's	A(A')=0	!A !B !C + !ABC + ABC
inverse	A(A')=0	!AB(C+!C) + ABC
DeMorgan's	AB = A+B	!AB + ABC
		B(!A+A C) * Absorption
		B(!A+C)

COMMON GATES



$X11 = 0xFFFFFB13$
 $SW X11 0(KS) \rightarrow X5 = FF FF FB 13$
 $sh X11 2(KS) \rightarrow X5 = FB 13 FB 13$
 $16x12 1(KS) + 1b = \text{sign extend} \rightarrow X5 = 0xFFFFFFFF$
 $16U x13 2(KS) + 1b = 0 \text{ extend} \rightarrow X5 = 0 + 00000000$
 load byte signed
 strai: divide by 2, round down



$$\text{min clock period} = \text{critical path}$$

$$\text{max clock freq} = \frac{1}{\text{min clock period}}$$

$$\begin{aligned} \text{MAX } t_{\text{clk-to-q}} + t_{\text{shortest-comb}} &\geq t_{\text{hold}} \\ \text{MIN } t_{\text{clk-to-q}} + t_{\text{longest-comb}} + t_{\text{setup}} &\leq t_{\text{q-period}} \end{aligned}$$

longer hold = value has potential to change
 shorter period = values may not finish computing

TIMING

SETUP: time before clock edge that input must be stable

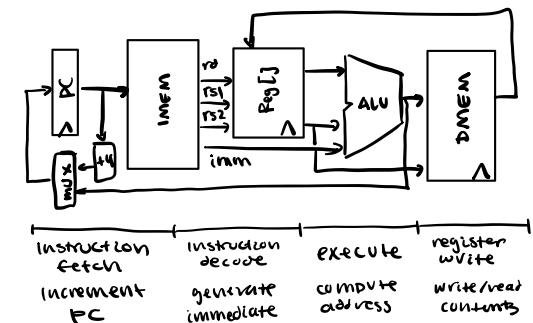
Hold: time after clock edge that input must be stable

CLK-to-Q: time after clock edge that input appears on Q

Critical path: longest path delay between 2 state elements

Path delay: $t_{\text{clk-to-q}} + \text{combinational logic}$ + $\frac{\text{setup time}}{\text{no element}}$

DATAPATH



$$O = PC + 4$$

$$I = \text{branch+jump}$$

IMMSet = type of immediate

ASel 0 = register
1 = PC (e.g. jalr) input for A ALU

BSel 0 = register
1 = immediate input for B ALU

ALUSet: 0=add
1=sll
2=slt
3=unused
4=xor
5=or
6=and
7=mul
8=mul
9=mul
10=unused
11=multiu
12=sub
13=sra
14=unused
15=bset

BRUn: 0=signed
1=unsigned
0:A=B
1:A=B

BRlt: 0=A>B
1=A<B
MEMRW: 0 = no write, but read
1 = write dataW @addr (stires S[b/h/w])

RegWEn: 0 = no writes back to Regfile allowed
1 = writes allowed

0 = mem data to register (loads I[b/bu/h/hu/w])

WBsel: 1 = ALU output to register (arithmetic)
2 = PC+4 to register (ra = destination register)

STRUCTURAL HAZARD: same piece of hardware being accessed by more instructions than access points

↳ SOLNS: stall / add more hardware, separate MEM & DMEM [can it be solved by adding another regfile?]

DATA HAZARD: data dependency → data needed by following instr before its updated ex reading register before prev finished writing

↳ SOLN: stall / NOP, forwarding

CONTROL HAZARD: flow of execution has more than 1 option

↳ SOLNS: JUMPS + BRANCHES

↳ SOLN: stall, branch prediction, flushing

INSTR	BREQ	BRlt	PCSel	ImmSel	BRUn	ASel	BSel	ALUSet	MemRW	RegWEn	WBsel
add	✓	✓	+4 ⁰	+	✓	Reg ⁰	Reg ⁰	Add	Read ⁰	1	ALU ¹
sub	✓	✓	+4	✓	✓	Reg	Reg	Sub	Read	1	ALU
addi	✓	✓	+4	I	✓	Reg	Imm ¹	Add	Read	1	ALU
lw	✓	✓	+4	I	✓	Reg	Imm	Add	Write ¹	0	ALU
sw	✓	✓	+4	S	✓	Reg	Imm	Add	Write	0	ALU
beq	0	✓	+4	B	✓	PC ¹	Imm	Add	Read	0	ALU
beq	1	✓	ALU ¹	B	✓	PC	Imm	Add	Read	0	ALU
bne	0	✓	ALU	B	✓	PC	Imm	Add	Read	0	ALU
bne	1	✓	+4	B	✓	PC	Imm	Add	Read	0	ALU
blt	✓	✓	ALU	B	signed 0 unsigned	PC	Imm	Add	Read	0	ALU
bltu	✓	✓	ALU	B	✓	PC	Imm	Add	Read	0	ALU
jalr	✓	✓	ALU	I	✓	Reg	Imm	Add	Read	1	
jal	✓	✓	ALU	J	✓	PC	Imm	Add	Read	1	
auipc	✓	-	+4	U	✓	PC	Imm	Add	Read	1	ALU

* longest instruction = load instructions

PIPELINING

MAX efficiency in n-stage pipeline = n instructions in pipeline

* add register for every wire carrying value across stages

latency = time for CPU to execute 1 instr

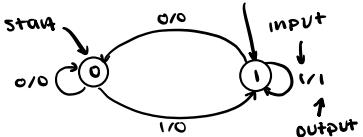
throughput = #of instructions executed per unit time

processor performance

$$\frac{\text{Time}}{\text{Program}} = \frac{\text{Instruction}}{\text{Program}} \cdot \frac{\text{Cycles}}{\text{Instruction}} \cdot \frac{\text{Time}}{\text{Cycles}}$$

	SINGLE	N-STAGE
CLOCK CYCLE	full datapath clock per	max of n stages
latency	clock cycle	clock cycle * n
throughput	$\frac{1}{\text{clock per full datapath}}$	$\frac{1}{\text{clock per max of n stages}}$

FSM = represented by states + transition
transition $i \rightarrow i$ output 1



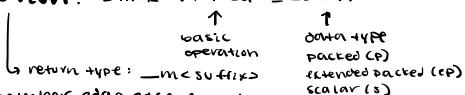
DLP = data level parallelism

* optimize instruction program → increase work per instruction

SIMD = single instr, multiple data

vectorized calculations

notation: - mm_ <instr_ops> <suffix>



↳ remember edge case for elements that don't fit in vector

AMDAHL'S LAW

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

F = fraction of program that can be optimized

S = speedup factor for how much that fraction can be optimized

OPENMP

#pragma omp parallel

```

{ 
  // fun stuff
  ...
  // copied & executed across all threads INDEPENDENTLY
  // if around a loop: copies + executes same loop across many threads
}
```

#pragma omp parallel for

```

for (l=0; l<h; l++) {
  // fun stuff
  ...
  // SPLIT across # of threads
  // each thread is independent
  // scheduled by OS
}
```

PLP : process level parallelism

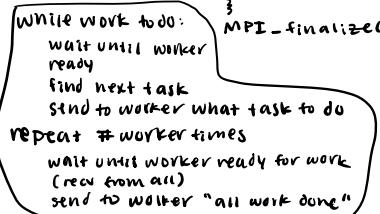
program on multiple processes @ once
Manager-worker framework

↳ manager assigns work every other process (workers) ask for work

OPENMPI

```

int main (int argc, char ** argv) {
  // set up openMPI
  if (manager process) {
    // manager node code
  } else {
    // worker node code
  }
  // terminate openMPI
}
```



CACHES

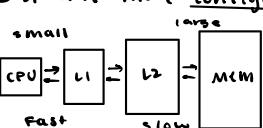
= small memory storage w/ FASTER ACCESS than DRAM

Logic:

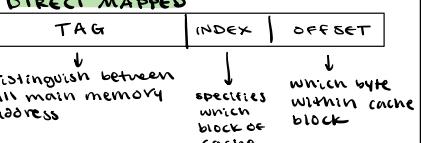
Need to access memory → check cache
if not in cache → grab from main mem & load into cache
if in cache → grab from cache

locality: principle for saving/updating data in cache

(1) temporal: keep most recently accessed items
(2) spatial: move contiguous blocks of memory into cache



DIRECT MAPPED



distinguish between all main memory address

of bits

$$\text{INDEX} = \log(\# \text{ blocks}/\text{associativity}) = \log_2(\#\text{of indices})$$

$$\text{offset} = \log_2(\text{size of blocks})$$

$$\text{tag} = (\#\text{bits in mem addr}) - (\#\text{of index bits}) - (\#\text{of offset bits})$$

$$\text{cache size} = \#\text{ of blocks} \times \text{block size}$$

example: summing vals ≥ 128

```

--m128i_l128 = _mm_set1_epi32(127);
--m128i_lsum_vec = _mm_setzero_si128();
for (int i=0; i < NUM_ELEMENTS/4 * 4; i+=4) {
  --m128i_ltmp = _mm_loadu_si128((__m128i *)vals+i); ← load 4 vals
  --m128i_lcomp_vec = _mm_cmplt_epi32(ltmp, _128i);
  --m128i_lmask = _mm_and_si128(ltmp, comp_vec); ← bitwise AND
  sum_vec = _mm_add_epi32(sum_vec, mask_vals);
}

int tmp_arr[4];
_mm_storeu_si128((__m128i *)tmp_arr, sum_vec);
result += tmp_arr[0] + tmp_arr[1] + tmp_arr[2] + tmp_arr[3];
for (int i=NUM_ELEMENTS/4 * 4; i<NUM_ELEMENTS; i+=4) {
  if (vals[i] >= 128) {
    result += vals[i];
  }
}

```

TLP = thread level parallelism

program = sequence of instrs

process = execution of program

↳ each w/ own address space + resources

↳ same code but **SEPARATE**

thread = segment of process

↳ same code, same time

↳ shared components: heap, global variables

INDIVIDUAL registers, PC, stack

↳ single core can execute **1** thread at a time

DATA RACE = multiple threads modifying same **SHARED VAR**

↳ soln: atomic operations + locks

LOCK: only 1 thread can have lock @ a time

- acquire: try to get lock + KEEP trying
- release: unlocks lock + continues
- try-acquire: try to get lock but if used → return FALSE

CRITICAL SECTIONS: code surrounded by lock → only ONE thread allowed to run

REDUCTIONS: increase parallelizability by reducing code in critical sections

t0 = 01111

s11i = 0111100000

srl1i = 0111100000000000

t0 = 0x F000 0000

t1 = 0x F000 0000

t0 = 0x DF00 0000

t1 = 0x FF00 0000

--m128i_l_mm_set1_epi32(first)

set 4 signed 32-bit ints to 1

--m128i_l_mm_loadu_si128(--m128i_l)

load into pointed to by p into vector

--m128i_l_mm_and_si128(--m128i_l, a, --m128i_l)

(a0, b0, a1, b1, a2, b2, a3, b3)

--m128i_l_mm_add_epi32(--m128i_l, a, --m128i_l)

(a0+b0, a1+b1, a2+b2, a3+b3)

void _mm_storeu_si128(--m128i_l, p)

store 128 vector at pointer p

--m128i_l_mm_and_si128(--m128i_l, a, --m128i_l)

bitwise AND of 128 bits in a and b

--m128i_l_mm_cmplt_epi32(--m128i_l, a, --m128i_l)

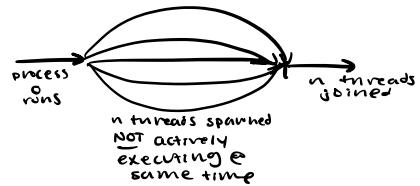
1st element of return vector set of 0xFFFFFFF

if 1st element of a and b = → 0x00000

↳ multithreaded **CANNOT** assume order threads are executed in

fork(): master → child threads & divides work

join(): wait until threads finish → synchronize → terminates all except master



$$kilo = 10^3 \text{ kibi} = 2^{10}$$

$$10^3 \approx 2^{10} \text{ mega} = 10^6 \text{ mibi} = 2^{20}$$

AMAT = average memory access time

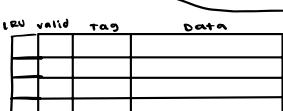
hit time + miss rate (miss penalty)

↓ time to get data if in cache already ↓ # misses ↓ time to get data if it's not in the cache

$$L2 \text{ AMAT} = L1 \text{ HT} + L1 \text{ MR} (L2 \text{ HT} + L2 \text{ MR} (\text{main mem access}))$$

cache code analysis

- ① count access per inner & outer loop
- ② look for compulsory misses, thrashing (going back and forth between filling cache block) and count hits + misses
- ③ Hit rate = $\frac{\text{hits}}{\text{hits} + \text{misses}}$



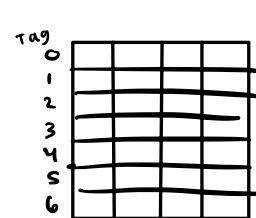
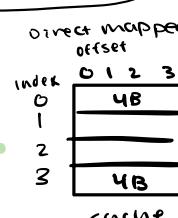
CACHE EVICTION POLICIES

LRU = least recently used

MRU = most recently used

FIFO = first in first out

LFU = last in first out



CACHE MISSES

① **compulsory miss**: never tried to pull block into cache (ALWAYS happens first time we request)

↳ SOLN: increase cache block size, prefetching

② **conflict miss**: had to evict block bc another block REPLACED it and CACHE NOT FULL

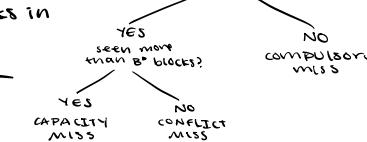
↳ does NOT exist for fully associative

↳ SOLN: increase associativity, improve replacement policy

③ **capacity miss**: evicted block bc cache was FULL

↳ increase cache size

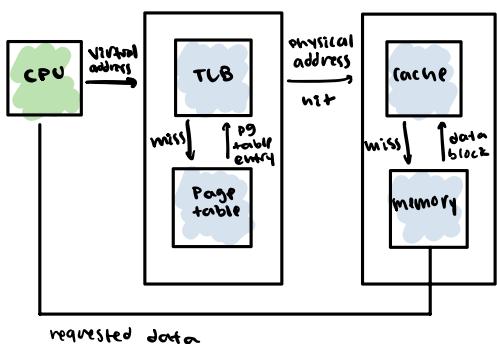
B = # total blocks in cache
= Cache size / block size



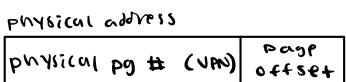
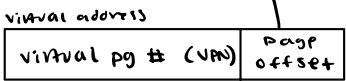
- cache = \$ main mem MM
- ① write through: on every write → data written to \$ + MM
 - ② no write allocate: \$ line NOT written to cache → directly to MM
 - ③ write back: data written to \$ and DIRTY BIT = 1 → written to MM in eviction
 - ④ write allocate: \$ line pulled to \$ and written to in specified write-hit policy

Virtual memory

- protection + security
- combine disk + memory into larger space
- allows programs to "see" more memory than exists
- PAGES** = uniformly sized block of data
- ↳ **page table** maps between virtual and physical



Offset bits = \log_2 (page byte size)
Div by 1 when write



TLB = translation lookaside buffer
 ↳ cache for page table
 ↳ usually fully associative
 ↳ RESET when new process started

0xB055ADE1
0xC160EFFF
0x7890ABCD
0x97543210
<truncated>
page table

0x00000445
 VPN offset

Index 0 = 0xB055ADE1

PPN + offset = ADE1445

MSB = 1 so valid

STEPS

① check if VPN in TLB
 ↳ check valid bit
 1 = hit
 0 = check page table

② check page table
 ↳ check valid bit
 1 = hit (NO page fault)
 0 = page fault
 ↓
 go to next physical address

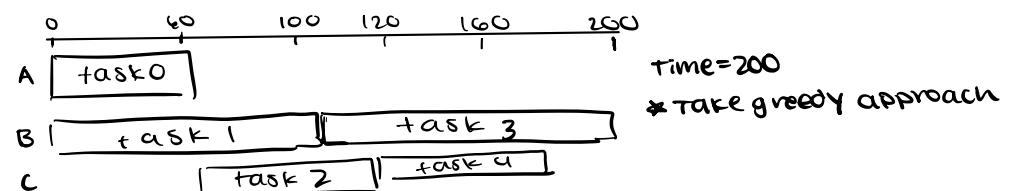
* FOR RISC + C coding questions :

- ③ can use bitwise AND (&) to check if a digit exists with another digit
- ④ can use slli to multiply
- ⑤ malloc(sizeof(int) * n)
- ⑥ calloc(num_elems, sizeof(int))

↑
 1 elem

MULTITHREADING TASKS

- * find combo of tasks to maximize parallelism (multiple chips doing work @ same time) and chip specialization (chip that's good @ memory doing tasks more demanding on memory)



Time = 200

* TAKE greedy approach