

PRINT your name: _____,
(last) (first)

PRINT your student ID: _____

You have 170 minutes. There are 10 questions of varying credit (100 points total).

Question:	1	2	3	4	5	6	7	8	9	10	Total
Points:	10	10	23	13	10	5	14	7	7	1	100

For questions with **circular bubbles**, you may select only one choice.

- Unselected option (completely unfilled)
- Only one selected option (completely filled)

For questions with **square checkboxes**, you may select one or more choices.

- You can select
- multiple squares
- (completely filled)

Anything you write that you ~~cross-out~~ will not be graded. Anything you write outside the answer boxes will not be graded. If you write multiple answers or your answer is ambiguous, we will grade the worst interpretation. For coding questions, you may write at most one statement and you may not use more blanks than provided.

If an answer requires hex input, make sure you only use capitalized letters! For example, `0xDEADBEEF` instead of `0xdeadbeef`. Please include hex (`0x`) or binary (`0b`) prefixes in your answers unless otherwise specified. For all other bases, do not add any prefixes or suffixes.

Read the following honor code and sign your name.

I understand that I may not collaborate with anyone else on this exam, or cheat in any way. I am aware of the Berkeley Campus Code of Student Conduct and acknowledge that academic misconduct will be reported to the Center for Student Conduct and may further result in, at minimum, negative points on the exam.

SIGN your name: _____

Q1 RISC-y Array Architecture**(10 points)**

Writing code to access integer arrays can be really annoying in RISC-V! Suppose we come up with new instructions, `readArr` to read from integer arrays and `writeArr` to write to integer arrays. For this question, you may assume integers are 32 bits.

`readArr rd, rs1, rs2` will read the array that `rs1` points to at the index stored in `rs2`, and put that value in register `rd`. In C pseudocode: `rd = ((int *) rs1)[rs2]`.

Q1.1 (3.5 points) What changes would we need to make to our datapath in order for us to implement the `readArr` instruction with as few changes as possible? Select all that apply.

- Add a new immediate type for ImmGen
- Add a new output to Regfile for a third register value
- Add a new input to the AMux and update any relevant selector/control logic
- Add a new input to the BMux and update any relevant selector/control logic
- Add a new ALU operation and update any relevant selector/control logic
- Add a new DMEM mux which feeds into the data input of the DMEM, and any relevant selector/control logic
- Add a new input to WBMux and update any relevant selector/control logic
- None of the above

(Question 1 continued...)

`writeArr rs3, rs1, rs2` will take the value in register `rs3`, and write that value to the array that `rs1` points to at index `rs2`. In C pseudocode: `((int *) rs1)[rs2] = rs3`.

Q1.2 (3.5 points) Assume that the changes, if any, for `readArr` have **not** been implemented for this subpart. What changes would we need to make to our datapath in order for us to implement the `writeArr` instruction with as few changes as possible? Select all that apply.

- Add a new immediate type for ImmGen
- Add a new output to Regfile for a third register value
- Add a new input to the AMux and update any relevant selector/control logic
- Add a new input to the BMux and update any relevant selector/control logic
- Add a new ALU operation and update any relevant selector/control logic
- Add a new DMEM mux which feeds into the data input of the DMEM, and any relevant selector/control logic
- Add a new input to WBMux and update any relevant selector/control logic
- None of the above

Q1.3 (3 points) Eddy noticed that the structure of `writeArr` is similar to an R-type instruction. However, when he tried to use the control signals for an R-type instruction, it didn't work. Which of the following control signals does he need to change to correctly implement `writeArr`? Select all that apply.

- PCSel
- ASel
- BSEL
- RegWEn
- MemRW
- None of the above

Q2 IF Only ID Pipelined Better**(10 points)**

In Project 3, we implemented a RISC-V CPU with two stages; stage 1 included IF and stage 2 included ID/EX/MEM/WB. For this question, imagine instead that we implement a two-stage pipeline with a different split; stage 1 will include IF/ID and stage 2 will include EX/MEM/WB (IF/ID/EX/MEM/WB are defined equivalently to the pipelined CPU on the reference card).

For Q2.1 and Q2.2, assume the following delays for each component. Any component not listed is assumed to have a negligible delay.

Component	Delay
$\tau_{\text{clk-to-q}}$	35ps
τ_{setup}	20ps
Mux	75ps
Regfile Setup	20ps
Regfile Read	175ps
Immediate Generator	150ps
Branch Comparator	200ps
ALU	200ps
Memory Read	300ps

Q2.1 (3 points) What is the minimum clock period of this circuit, in picoseconds, to achieve correct behavior?

ps

Q2.2 (2 points) Which component in stage 2 can we move to stage 1 to decrease the minimum clock period of this circuit the most, while maintaining the same behavior? If a decrease is not possible, write "Not Possible".

Component:

For the remainder of this question, assume that the changes made in Q2.2, if any, have **not** been implemented.

Q2.3 (3.5 points) In the CPU, which of the following values must have a pipeline register? Select all that apply.

- | | | |
|--|---------------------------------------|--|
| <input type="checkbox"/> Instruction | <input type="checkbox"/> RegReadData2 | <input type="checkbox"/> MemReadData |
| <input type="checkbox"/> Program Counter | <input type="checkbox"/> Immediate | <input type="checkbox"/> None of the above |
| <input type="checkbox"/> RegReadData1 | <input type="checkbox"/> ALUOut | |

Q2.4 (1.5 points) Assume that the pipeline has been correctly implemented. Which types of hazards could a program experience? Assume that you cannot read from and write to the Regfile in the same clock cycle.

- | | | | |
|----------------------------------|-------------------------------|-------------------------------------|-------------------------------|
| <input type="checkbox"/> Control | <input type="checkbox"/> Data | <input type="checkbox"/> Structural | <input type="checkbox"/> None |
|----------------------------------|-------------------------------|-------------------------------------|-------------------------------|

This page intentionally left (mostly) blank.

The exam continues on the next page.

Q3 Faster Than Average**(23 points)**

In this question, you will parallelize a function to compute the average of all values in a matrix. Below is a single-threaded implementation of this function.

```

1 double matrix_average(double** matrix, int num_rows, int num_cols) {
2     double global_sum = 0.0;
3     for (int i = 0; i < num_rows; i++) {
4         for (int j = 0; j < num_cols; j++) {
5             global_sum += matrix[i][j];
6         }
7     }
8     return global_sum / (num_rows * num_cols);
9 }

```

Using the SIMD operations provided, optimize `matrix_average`. You have access to the following SIMD operations. A vector is a 256-bit vector register capable of holding 4 doubles.

- `vector vec_load(double* A)`: Loads 4 doubles at memory address `A` into a vector
- `void vec_store(double* A, vector B)`: Stores the 4 doubles in vector `B` at memory address `A`
- `vector vec_set0()`: Puts all 0s into a vector
- `double vec_sum(vector A)`: Adds all elements of the vector together: `return A[0] + A[1] + A[2] + A[3]`
- `vector vec_add(vector A, vector B)`: Adds `A` and `B` together elementwise

```

1 double matrix_average(double** matrix, int num_rows, int num_cols) {
2     double global_sum = 0.0;
3     vector sum_vec = _____;
4     for (int i = 0; i < num_rows; i++) {
5         for (int j = 0; j < _____; _____) {
6             vector values = _____;
7             sum_vec = _____;
8         }
9         for (int j = _____; j < _____; _____) {
10            global_sum += matrix[i][j];
11        }
12    }
13    global_sum += _____;
14    return global_sum / (num_rows * num_cols);
15 }

```

(Question 3 continued...)

Parallelize `matrix_average` using OpenMP without using `#pragma omp parallel for` or `reduction`. Each thread should work on one or more rows of the matrix. Assume `num_rows` is a multiple of `num_threads`.

```
1 double matrix_average(double** matrix, int num_rows, int num_cols) {
2     double global_sum = 0.0;
3
4     _____ Q3.10
5     {
6         int num_threads = omp_get_num_threads();
7         int thread_num = omp_get_thread_num();
8
9         int chunk_size = _____; Q3.11
10
11        int start_row = _____; Q3.12
12
13        int end_row = _____; Q3.13
14
15        for (int i = start_row; i < end_row; i++) {
16            double row_sum = 0.0;
17            for (int j = 0; j < num_cols; j++) {
18                row_sum += matrix[i][j];
19            }
20            _____ Q3.14
21            _____ Q3.15
22        }
23    }
24    return global_sum / (num_rows * num_cols);
25 }
```

Q4 *Convoluting Caching*

(13 points)

Consider the following function that takes in two integer arrays, *a* (of length *a_len*) and *b* (of length *b_len*), and returns the 1D convolution of *a* and *b*. Assume *results* is properly allocated.

Let *a*=0x1000, *b*=0x2000, *results*=0x3030, *a_len*=4, and *b_len*=2.

```
void convolve_1d(int* a, int a_len, int* b, int b_len, int* results) {
    for (int i = 0; i < a_len - b_len + 1; i++) {
        register int sum = 0;
        for (int j = 0; j < b_len; j++) {
            sum += b[j] * a[i + j];
        }
        results[i] = sum;
    }
}
```

For Q4.1 and Q4.2, we have a single-level, *direct-mapped* 64B cache with 16B blocks and 16-bit addresses.

Q4.1 (3 points) What are the tag, index, and offset bits of the address 0x3037?

T: 0b

I: 0b

O: 0b

Q4.2 (2.5 points) What is the overall hit rate for a call to *convolve_1d*? No need to simplify the fraction.

Q4.3 (2.5 points) We change to a *2-way set associative* cache of the same size with a LRU replacement policy. What is the overall hit rate for a call to *convolve_1d*? No need to simplify the fraction.

Q4.4 (2.5 points) We change to a *fully associative* cache of the same size with a LRU replacement policy. What is the overall hit rate for a call to *convolve_1d*? No need to simplify the fraction.

Q4.5 (2.5 points) We discover that accessing physical memory will take 400 cycles, so we decide to add an L2 cache. The hit rate of the L1 cache is 75%, and the hit rate of the L2 cache is 99%. With an access time of 6 cycles to fetch from the L1 cache, and an access time of 36 cycles to fetch from the L2 cache, what would our memory access time be for this system, on average?

cycles

Q5 The Lookup Box**(10 points)**

Consider a system with a 32-bit virtual address space, 256B pages, and 16 MiB of DRAM as main memory. Provided below is the TLB and a portion of the page table. The TLB is fully associative and there is no data cache. The next free physical pages in main memory start at physical addresses 0x61DE00 and 0x61EF00, respectively.

Each PTE is 32 bits. Bit 31 is the valid bit, bit 30 is the dirty bit, bits 16 through 29 hold other metadata (not relevant for this question), and bits 0 through 15 hold the PPN.

Tag (VPN)	PPN	Valid	Dirty
0x000000	0x23EF	1	0
0x000001	0xFFFF	0	0

Index	PTE
0x0	0x80AB23EF
0x1	0x80EE00C0
0x2	0x8123200A
0x3	0x3561CBA8
...	...
0xA	0xCAFFEEEE

For each question, determine what the memory address access results in, and calculate its physical address. Note that each memory access is executed in sequence, so they are **not** independent of each other.

Q5.1 (2.5 points) Virtual Address: 0x000000FF

- TLB hit
 TLB miss, no page fault
 TLB miss, page fault

Physical Address:

0x

Q5.2 (2.5 points) Virtual Address: 0x00000283

- TLB hit
 TLB miss, no page fault
 TLB miss, page fault

Physical Address:

0x

Q5.3 (2.5 points) Virtual Address: 0x00000AAA

- TLB hit
 TLB miss, no page fault
 TLB miss, page fault

Physical Address:

0x

Q5.4 (2.5 points) Virtual Address: 0x00000360

- TLB hit
 TLB miss, no page fault
 TLB miss, page fault

Physical Address:

0x

Q6 *Cumulative: Potpourri*

(5 points)

Q6.1 (1 point) The OS running on a cluster of computers in a datacenter allows a single machine to read and write the memory and local disk of a remote machine in the same rack or array.

According to lecture, which of the following are true? Select all that apply.

("farther away" means that the distance the data travels increases in steps, first to our local machine, then to a machine in our same rack, then to a machine in our same array.)

- As our CPU sends data to DRAM farther away, bandwidth increases
- As our CPU sends data to Disk farther away, bandwidth increases
- As our CPU sends data to DRAM farther away, latency increases
- As our CPU sends data to Disk farther away, latency increases
- We have higher latency to DRAM on an Array computer than to our own disk
- We have higher bandwidth to DRAM on an Array computer than to our own disk
- None of the above

Q6.2 (1 point) After a single machine M finishes its assigned portion of a map task in a MapReduce cluster, which of the following can happen immediately, regardless of the overall program state? Select all that apply.

- M can shut down without risking the success of the overall computation
- M can be assigned a new map task
- Data shuffling of the workload M just finished can begin
- The reduce task for the workload M just finished can begin
- None of the above

Q6.3 (1 point) We want to send **one** bit using a Hamming error correcting code. What are the valid bit patterns you could send that correspond to 0b0 and 0b1?

0b0: 0b1:

(Question 6 continued...)

Your network card just received a packet with an incorrect checksum.

Q6.4 (1 point) According to lecture, which of the following is true?

- There was a guaranteed error in the payload but not the checksum
- There was a guaranteed error in the checksum but not the payload
- There was a guaranteed error in either the payload or checksum
- None of the above

Q6.5 (1 point) According to lecture, what should you do?

- Send back a traditional data packet with information about which packet had the problem
- Send back an "ACK"
- Send back a "NO-ACK"
- Send back nothing and delete the packet

Q7 Cumulative: RV32tok**(14 points)**

Write a program `splitCode`, which will split a RISC-V program into blocks of code with no branches or jumps (jal or jalr). Specifically, `splitCode` will have the following function signature:

- **Input:** `int* code`, an array of RISC-V instructions. Each RISC-V instruction is stored as a 32-bit integer, equal to its translation. You may assume that all instructions are valid RISC-V base instructions, and that there are no pseudoinstructions, `ecalls`, or `ebreaks`.
- **Input:** `n`, the number of instructions in `code`.
- **Input:** `int*** result`, a pointer to store your result. Your result should be an array of `int*`s, where each `int*` points to the beginning of a sequence of consecutive instructions with no branches or jumps. Each of these arrays should be "null-terminated"; that is, the last element of each array should be the number 0, to signify the end of the array. Every non-branch/non-jump instruction must be represented in exactly one subarray of your result. No branch/jump instruction should be in any subarray of your result.
- **Output:** `int`, the length of your result.

For example, for the following RISC-V code:

```
1 beq x0 x0 pass
2 beq x0 x0 pass
3 add a0 t0 t1
4 add t0 a0 a1
5 add t0 a1 a2
6 xor a0 t0 t1
7 j pass
8 addi t0 x0 1
9 addi t0 x0 2
10 beq x0 x0 pass
```

`result` should point to the following array, and the return value should be 5.

```
[
  // The instructions before line 1
  [0],
  // The instructions between lines 1 and 2
  [0],
  // The instructions between lines 2 and 7
  [add a0 t0 t1, add t0 a0 a1, add t0 a1 a2, xor a0 t0 t1, 0],
  // The instructions between lines 7 and 10
  [addi t0 x0 1, addi t0 x0 2, 0],
  // The instructions after line 10
  [0]
]
```

Useful C function prototypes:

```
void* malloc(size_t size);
void* calloc(size_t num_elements, size_t size);
void* memcpy(void* dest, void* source, size_t num_bytes);
```

(Question 7 continued...)

```
1 // Returns true if instruction is a branch or jump instruction
2 bool isBranchJump(int instruction) {
3     return _____;
4 }
5
6 int splitCode(int* code, int n, int*** result) {
7     int num = 0; // total number of branches and jumps
8     for(int i = 0; _____; i++) {
9         num += _____;
10    }
11    int** data = malloc(_____);
12    int* codecopy = calloc(n+1, _____);
13    // Hint: You should not need any more memory allocations
14    memcpy(codecopy, code, _____);
15    for(int i = 0; _____; i++) {
16        data[i] = _____;
17        while(_____ && _____ != 0) {
18            _____;
19        }
20        _____;
21        codecopy++;
22    }
23    _____;
24    return _____;
25 }
```

Q8 Cumulative: Chips Ahoy**(7 points)**

Consider the following set of tasks:

Task ID	Time (minutes)	Prerequisites	Time Breakdown
0	100	-	90% memory, 10% math
1	100	-	90% memory, 10% math
2	100	0	30% memory, 70% math
3	100	0,1	30% memory, 70% math
4	100	1	30% memory, 70% math
Total	500	-	54% memory, 46% math

After running this set of tasks on your local, single-threaded CPU (referred to as Chip A) in 500 minutes, you decide that it's too slow and decide to upgrade to a new chip.

At the store, you find two options:

- Chip B: A single-threaded CPU optimized for memory accesses. It can do memory operations 3 times as fast as Chip A, but it takes twice as long to do math operations.
- Chip C: A single-threaded GPU optimized for fast math. It can do math operations practically instantly (infinite times speedup), but it takes twice as long as Chip A to do memory operations.

Q8.1 (2 points) Using Chip B alone, how many minutes would all 5 tasks take?

Q8.2 (2 points) Using Chip C alone, how many minutes would all 5 tasks take?

Q8.3 (3 points) Using one chip was still taking too long, so you buy both Chip B and Chip C from the store, and connect them to Chip A in a new multicore machine with negligible overhead. Using all three chips, what is the minimum amount of time required to complete this set of tasks? Each task must be completed entirely on one chip.

Please also provide the list of tasks each chip will complete, in order of completion, or write "None" if a chip does not complete any task. For example, if you decide to have Chip A complete all of the tasks, your answer should be "0, 1, 2, 3, 4" for Chip A, and "None" for Chip B and Chip C.

Q9 *Cumulative: The Magnus Effect*

(7 points)

Q9.1 (7 points) Write a Boolean expression that determines if a 19-bit unsigned integer can be expressed exactly as a 19-bit floating point number. For full credit, you may use at most 8 Boolean operators ($|$, $\&$, \sim).

Inputs: Bits A through S

Output: One bit. Output 1 if $0b\ ABC\ DEFG\ HIJK\ LMNO\ PQRS$ is an unsigned number that can be exactly represented as a 19-bit float which follows all IEEE-754 conventions, with 5 exponent bits (and a standard bias of -15). Output 0 otherwise.

Hint: both the exponent and the mantissa provide nontrivial constraints.

For partial credit, describe in English how you can determine if a 19-bit unsigned integer can be expressed exactly as a 19-bit floating point number (using the float representation described above).

Q10 *The Finish Line*

(1 points)

Everyone will receive credit for this question, even if you leave it blank.

Q10.1 (1 point) On a scale of 1 to 10, rate Eddy's weekly announcement puns.

Q10.2 (0 points) Is there anything you want us to know?