

ECE 454 – Computer Systems Programming

The Edward S. Rogers Sr. Department of Electrical and Computer Engineering

Midterm Examination Fall 2010

Name	
Student #	

Answer all questions. Write your answers on the exam paper. Show your work.
Each question has a different assigned value, as indicated.

The exam is open book (only simple calculators allowed, no cell phones or PDAs)

Total time available: **110 minutes**

Total marks available: **100** (roughly one mark per minute, with 10 extra minutes)

Verify that your exam has all the pages.

Part	Points	Mark
1	8	
2	8	
3	14	
4	10	
5	10	
6	20	
7	20	
8	10	
Total	100	

PART 1) [8] Measuring Programs

Which measuring tool should you use in each of the following scenarios? For each choose the simplest and lowest overhead tool that is effective for the task. Put the number for the tool in the box next to each scenario.

List of Tools:

- 1) command-line timer like `/usr/bin/time`
- 2) C library timer like in `<sys/times.h>`
- 3) Hardware timer or performance counters
- 4) General-purpose instrumentor like PIN
- 5) GPROF/GCOV
- 6) Simulator (hardware or software)

Scenarios: for each you want to measure...

- | | |
|---|--|
| a) The minimum number of processor cycles to execute a function | <div style="border: 1px solid black; padding: 2px; display: inline-block; text-align: center;">3</div> |
| b) The average duration of a long-running program | <div style="border: 1px solid black; padding: 2px; display: inline-block; text-align: center;">1</div> |
| c) The number of instructions per-instance of a certain function (on average) | <div style="border: 1px solid black; padding: 2px; display: inline-block; text-align: center;">4</div> |
| d) The impact of a new coherence scheme on performance of a program | <div style="border: 1px solid black; padding: 2px; display: inline-block; text-align: center;">6</div> |
| e) The top function (by time) of a long-running program | <div style="border: 1px solid black; padding: 2px; display: inline-block; text-align: center;">5</div> |
| f) The duration of a function within a 1-second program | <div style="border: 1px solid black; padding: 2px; display: inline-block; text-align: center;">2</div> |
| g) The cache miss behavior of a production web-server in-the-field (i.e., live) | <div style="border: 1px solid black; padding: 2px; display: inline-block; text-align: center;">3</div> |
| h) The average behavior of an if-else branch for a hot function in a long-running program | <div style="border: 1px solid black; padding: 2px; display: inline-block; text-align: center;">5</div> |

PART 2) [8] Short Answer

- a) Name six things that a modern superscalar processor (such as the ones in the 'UG' machines) does to optimize the execution of a program automatically and transparently to the programmer. Just the one or two word name for each is fine.

3marks 0.5 each

Cache
Instruction Scheduling (ooo issue, early loads)
Prefetching
Parallel/wide issue/execution
Pipelining
Branch prediction
TLB

- b) What is the difference between coherence and consistency? A one or two sentence answer should suffice.

2 marks

Coherence is with respect to a single location while consistency is with respect to multiple locations. Both refer to the rules/method by which memory locations are managed by a parallel machine.

- c) Give three reasons that you hypothetically might NOT want to use -O3 optimization for gcc when compiling a program?

3 marks 1each

You care about compilation time
You care about program size
-O3 invokes a compiler bug or a bug in your program
You are debugging

PART 3) [14] Optimization Decisions

The following shows the (fake) top four results of a GPROF measurement of VPR. OptsRus has three optimizations planned, and want to maximize the performance of VPR.

% time	cumulative seconds	self seconds	calls	self s/call	total s/call	name
66.00	1.71	1.71	891388	0.00	0.00	try_swap
16.00	2.13	0.42	26666027	0.00	0.00	comp_td_p2p_delay
5.00	2.25	0.12	68548	0.00	0.00	label_wire_muxes
3.00	2.33	0.08	3160610	0.00	0.00	update_bb

Optimization1: will make try_swap 10% faster (i.e., 1.1x faster)

Optimization2: will make comp_td_p2p_delay 100% faster (i.e., 2.0x faster)

Optimization3: will inline comp_td_p2p_delay into update_bb, but the new update_bb will be 200% slower (i.e., 3x slower).

- a) OptsRus only has time/resources to implement one optimization. Which optimization should you choose to implement and why? Give the final expected program speedup of your chosen optimization. HINT: your answer should involve some simple calculations, and should give the results of those.

OptsRus should implement Optimization # 3

Because:

Opt1: $66/1.1 + 34 = 94$ % of original time (1.06)

Opt2: $16/2.0 + 84 = 92$ % of original time (1.08)

Opt3: $3*3 + 81 = 90$ % of original time (1.11)

b) If OptsRus found the time/resources to implement a second optimization, which should it be and why?

OptsRus should implement Optimization # 1

Because:

Optimization 2 inlines comp_td_p2p_delay, so it is unclear you can still make it go faster after inlining. But we also accept the answer of choosing optimization #2, if you do reasonable calculations showing that it improves speedup better than optimization #1 under clear assumptions.

PART 4) [10] Compiler Optimization

Name 8 optimizations that you would hope your compiler would do to this code. For each, give the formal name of the optimization and very briefly describe which variable(s) or code parts it will modify/improve and how. The first answer of 8 is given as an example.

```
1: x=5;
2: y=2;
3: debug = 0;
4: z=x+y;
5:
6: for (i=0;i<100;i++){
7:     m += i*x+y;
8:     n = z*y;
9:     A[i] += A[m];
10:
11:     if (debug){
12:         printf("m: %d\n",m);
13:     }
14:     q += i*x;
15: }
16:
17: printf("result: %d %d %d %d %d %d\n",m,n,q,x,y,z);
```

constant propagation: line4 changed to z=5+2

- 1 assign m,n,q to registers to reduce load/store latency
- 1 constant folding: line4 changed to z=7
- 1 constant propagation: line7 changed to m+= i*5+2
- 1 constant propagation: line8 changed to n = 7*2
- 1 constant folding: line8 changed to n = 14
- 3 loop invariant code motion: line8 moved outside of loop
- 1 dead code elimination: lines 11-13 deleted (cannot execute since debug=0)
- 2 common subexpression elimination: both line7 and line14 compute i*x (i*5 after ConstProp)
- 1 line9 is dead (if you assume that A[] is unused later in the code)
- 1 const propagation line14 changes to i*5
- 1 strength reduction: line8 becomes n = z<<1 (if you didn't propagate a constant for z)
- 1 dead code elimination: line3 is dead if you eliminated lines11-13
- 1 constant propagation: line11 if(0)
- 1 reorder instructions to hide the latency of loading A[m], or prefetch A[m], or q, or A[i]

NOTE: not all of these make sense together; and of course the most you can get is 10 points

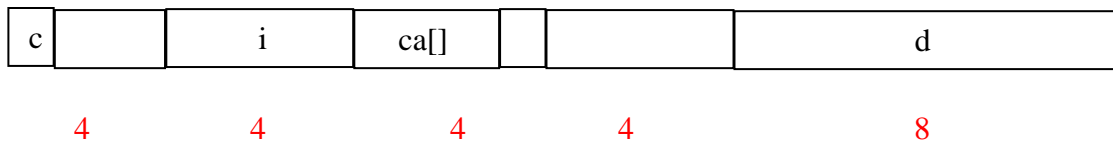
PART 5) [10] Static Memory

Consider the following declaration, Assume a CPU architecture like the UG machines, with typical 32-bit x86 linux.

```
struct S {  
    char c;  
    int i;  
    char ca[3];  
    double d;  
} A[100];
```

- a) How much space (in bytes) does the following array of structs take in memory? Draw (horizontally) the layout of the elements.

5marks:

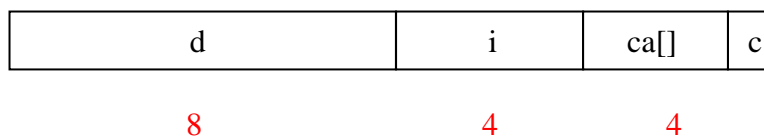


i must be 4-byte-aligned, d must be 8-byte-aligned
24B per struct, 2400B total for array.

- b) Can you modify the declaration to save space? If so, give the new declaration, and give the size of the modified array. Draw (horizontally) the new layout of the elements.

5marks:

```
struct S {  
    double d;  
    int i;  
    char ca[3];  
    char c;  
} A[100];
```



16B per struct, 1600B total

PART 6) [20] Memory Performance

OptsRus must optimize the following variant of matrix multiply code to target a certain machine, the specs for which are below.

```
// assume A[N][N], B[N][N], C[N][N], and D[N][N] have integer elements

for (i=0;i<N;i++){
    for (j=0;j<N;j++){
        for (k=0;k<N;k++){
            A[i][j] += B[i][k] * C[k][j] + D[i][j];
        }
    }
}
```

Target Machine Specs:

16KB L1 Cache

4MB L2 Cache

64B cache blocks (for both caches)

From past experience, OptsRus knows that tiling is the way to go. Considering tiling as the *only* optimization, and also considering only square tiles, what should the tile dimension T (assuming TxT tiles) be to optimize memory performance for each of the given values of N (array dimension) for the code? Give the value of T as both the number of elements and the number of bytes. Show the work you did to compute the tile dimension T and explain in one sentence why this choice should work well. You do NOT have to show the tiled code.

NOTE: we are assuming that each tile size you choose would be theoretically the best, ignoring prefetching and other subtleties of the machine that might lead to a different answer via experimentation as in the lab homework.

a) $N = 512$ T : _____elements, _____bytes

Array size = $(512\text{elements})^2 * 4\text{B/element} = (2^{18}) * (2^2) = 2^{20} = 1\text{MB}$ (per array)

4 Arrays = $4 * 1\text{MB} = 4\text{MB}$ total for all four arrays

This fits exactly in the L2, hence there will be no L2 capacity misses, hence we should tile for the L1:

Tile size = $L1\text{-size}/4\text{arrays} = 16\text{KB}/4\text{arrays} = 4\text{KB/array}$

$4\text{KB}/(4\text{B/element}) = 1024$ elements

Tile dimension $T = \sqrt{1024} = \sqrt{2^{10}} = 2^5 = 32\text{elements}$

Therefore $T = 32$ elements $32\text{elements} * (4\text{B/element}) = 128\text{B}$

b) $N = 1024$ T : _____elements, _____bytes

Array size = $(1024\text{elements})^2 * 4\text{B/element} = (2^{20}) * (2^2) = 2^{22} = 4\text{MB}$ (per array)

4 Arrays = $4 * 4\text{MB} = 16\text{MB}$ total for all four arrays

This is much bigger than the L2, so we should tile for the L2.

Tile size = $L2\text{-size}/4\text{arrays} = 4\text{MB}/4\text{arrays} = 1\text{MB/array}$

$1\text{MB}/(4\text{B/element}) = 2^{20}/(2^2) = 2^{18}\text{elements}$

Tile dimension $T = \sqrt{2^{18}} = 2^9 = 512$ elements

Therefore $T = 512$ elements $512\text{elements} * (4\text{B/element}) = 2048\text{B}$

PART 7) [20] Dynamic Memory Allocation

Consider an allocator with the following specification:

- Uses a single explicit free list.
- All memory blocks have a size that is a multiple of 8 bytes and is at least 16 bytes.
- All headers, footers, and pointers are 4 bytes in size
- Allocated blocks consist of a header and a payload (no footer)
- Free blocks consist of a header, two pointers for the next and previous free blocks in the free list, and a footer at the end of the block.
- All freed blocks are immediately coalesced.
- The heap starts with 0 bytes, never shrinks, and only grows large enough to satisfy memory requests.
- The heap contains only allocated blocks and free blocks. There is no space used for other data or special blocks to mark the beginning and end of the heap.
- When a block is split, the lower part of the block becomes the allocated part and the upper part becomes the new free block.
- Any newly created free block (whether it comes from a call to free, the upper part of a split block, or the coalescing of several free blocks) is inserted at the beginning of the free list.
- All searches for free blocks start at the head of the list and walk through the list in order (i.e., first-fit).
- If a request can be fulfilled by using a free block, that free block is used. Otherwise the heap is extended only enough to fulfill the request. If there is a free block at the end of the heap, this can be used along with the new heap space to fulfill the request.

Below you are given a series of memory requests. You are asked to show what the heap looks like after each request is completed using a first fit placement policy. The heap is represented as a series of boxes, where each box is a single block on the heap, and the bottom of the heap is the left most box. In each block, you should write the total size (including headers and footers) of the block in bytes and either 'f' or 'a' to mark it as free or allocated, respectively. For example, the following heap contains an allocated block of size 16, followed by a free block of size 32.

16a	32f
-----	-----

Assume that the heap is empty before each of the sequences is run, with a single 200B free block (shown for you). You do not necessarily have to use all the boxes provided for the heap. Some of the boxes are already filled in to help you. It is recommended to solve this on a scrap paper then copy your final answer into the boxes when you are satisfied. There are two copies here in case you make a mess. Clearly circle the one you want graded and cross out the one you don't want graded.

COPY #1: (they are identical, we will only grade the one you circle)

	200f					
ptr1 = malloc(20)	24a	176f				
ptr2 = malloc(30)	24a	40a	136f			
free(ptr1)	24f	40a	136f			
free(ptr2)	200f					
ptr3 = malloc(72)	80a	120f				
ptr4 = malloc(1)	80a	16a	104f			
ptr5 = malloc(12)	80a	16a	16a	88f		
ptr6 = malloc(24)	80a	16a	16a	32a	56f	
free(ptr3)	80f	16a	16a	32a	56f	
free(ptr5)	80f	16a	16f	32a	56f	
ptr7 = malloc(8)	80f	16a	16a	32a	56f	

NOTE: since there is an explicit free list, the first-fit free block is not necessarily the first free block in the left-to-right order!

COPY #2: (they are identical, we will only grade the one you circle)

	200f					
ptr1 = malloc(20)						
ptr2 = malloc(30)						
free(ptr1)						
free(ptr2)						
ptr3 = malloc(72)						
ptr4 = malloc(1)						
ptr5 = malloc(12)						
ptr6 = malloc(24)						
free(ptr3)						
free(ptr5)						
ptr7 = malloc(8)						

PART 8) [10] Locks and Critical Sections

For the following code, explain what you could do to improve its performance without changing its output (i.e., so it still works).

Assumptions:

- These are the only three threads in the program, running in parallel on separate CPUs.
- v,w,x,y,z are all in shared memory
- You are not allowed to add additional critical sections.
- You do not know what is inside the ‘update’ functions, you cannot optimize them.

HINT: think about what you learned about locks, sharing, and coherence behavior.

```
int v;  
int w;  
int x;  
int y;  
int z;
```

Thread1:

```
while(1){  
    ...  
    test_&set_lock(L);  
    y = updatey(y);  
    w = updatew(w);  
    z = updatez(z);  
    unlock(L);  
}
```

Thread2:

```
while(1){  
    ...  
    test_&set_lock(L);  
    y = updatey(y);  
    v = updatev(v);  
    w = updatew(w);  
    z = updatez(z);  
    x = updatex(x);  
    unlock(L);  
}
```

Thread3:

```
while(1){  
    ...  
    test_&set_lock(L);  
    y = updatey(y);  
    x = updatex(x);  
    z = updatez(z);  
    unlock(L);  
}
```

Use a better lock, at least a test&test&set with backoff or better

v= updatew(v); can be moved out of the critical section for Thread2 (only one updating v)

V,w,x should be padded to all be on separate cache blocks since they are updated separately

Y and z should stay on the same cache block since they are always updated together.