

University of Pennsylvania
Department of Electrical and System Engineering
System-on-a-Chip Architecture

ESE532, Fall 2018

Final

Friday, December 14

- Exam ends at 11:00AM; begin as instructed (target 9:00AM).
Do **not** open exam until instructed.
- Problems weighted as shown.
- Calculators allowed.
- Closed book = No text or notes allowed.
- Show work for partial credit consideration.
- Unless otherwise noted, answers to two significant figures are sufficient.
- Sign Code of Academic Integrity statement (see last page for code).

I certify that I have complied with the University of Pennsylvania's Code of Academic Integrity in completing this exam.

Name: [Solution](#)

1a	1bc	2a	2b	2c	3	4	5	6a	6b	6c
5	5	2	2	8	9	10	9	3	3	4

7a	7b	7cd	7e	8a	8b	8c	8d	Total
6	2	6	8	6	2	6	4	100

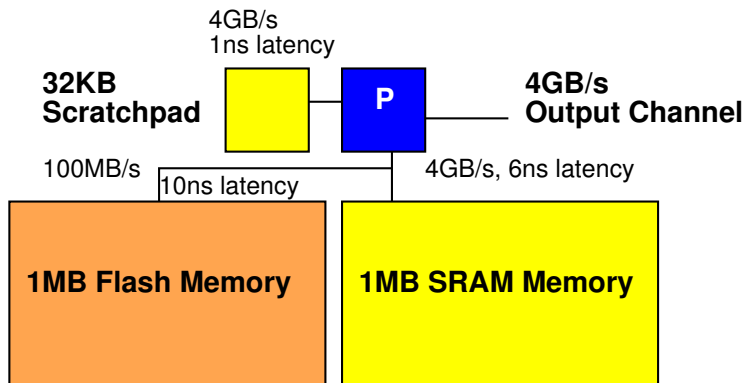
[Average 58, Std. Dev. 11](#)

```

// You will be determining a value for FREQBYTES
#define WINDOW 1024
#define MAXBITLEN 11
#define LOG_MAXBITLEN 4
#define MAX_FREQS 255
#define MASKLOOKUP ((1<<MAXBITLEN)-1)
#define MASKLEN ((1<<LOG_MAXBITLEN)-1)
#define AMPLEN 14
#define FREQLEN 14
#define MASKAMP ((1<<AMPLEN)-1)
#define MASKFREQ ((1<<FREQLEN)-1)
uint8_t in[FREQBYTES];
uint32_t fa[FREQS];
uint32_t lookup[1<<MAXBITLEN];
uint16_t s[MAX_FREQS][WINDOW];
while(1) { // Outer while loop
    uint32_t ts[WINDOW];
    for (j=0;j<WINDOW;j++) ts[j]=0; // Loop A
    uint8_t freqs=read_flash_byte(); // max rate 100MB/s
    for(int i=0;i<FREQBYTES;i++) // Loop B
        in[i]=read_flash_byte();
    uint11_t top11=((int *)in)[0]>>21;
    uint11_t next11=(((int *)in)[0]>>10)&MASKLOOKUP;
    int next11bitpos=11;
    for(i=0;i<freqs;i++) { // freqs<MAX_FREQS // Loop C
        uint32_t res=lookup[top11];
        uint32_t tfa=res>>LOG_MAXBITLEN; fa[i]=tfa;
        uint4_t len=MASKLEN & res;
        uint32_t t1=(top11<<len); uint4_t t2=(MAXBITLEN-len); uint32_t t3=(next11>>t2);
        top11= t1|t3;
        next11bitpos+=len;
        uint32_t bytupos=next11bitpos>>3; uint3_t bitoffset=next11bitpos%8;
        uint32_t wordval=*((int *)&in[bytupos])); // treat as 1 cycle
        uint4_t t4=(21-bitoffset); uint32_t t5=(wordval>>t4);
        next11=MASKLOOKUP & t5;
    }
    for (i=0;i<freqs;i++) { // Loop D
        uint16_t freq=(fa[i]>>AMPLEN) & MASKFREQ;
        uint16_t amp=fa[i] & MASKAMP;
        for (j=0;j<WINDOW;j++) // Loop E
            ts[j]+=s[freq][j]*amp;
    }
    for (j=0;j<WINDOW;j++) // Loop F
        output(ts[j]); // max rate 4GB/s
}

```

We start with a baseline, single processor system as shown.



- For simplicity throughout, we will treat non-memory indexing adds (subtracts count as adds), shifts, mod-by-power-of-two, ORs, ANDs, and multiplies as the only compute operations. We'll assume the other operations take negligible time or can be run in parallel (ILP) with the adds, multiplies, and memory operations. (Some consequences: You may ignore loop and conditional overheads in processor runtime estimates; you may ignore computations in array indices.)
- Assume all additions are associative.
- Baseline processor can execute one compute operation (above) per cycle and runs at 1 GHz.
- Constant expressions (like $1 \ll 8$) are evaluated by the compiler and take no time to compute at runtime.
- Maximum data rate for reading from flash is 100MB/s. Latency of read is 10 ns.
- The output port used by `output()` can transfer data at 4GB/s (one 32b word per cycle at 1 GHz).
- Baseline processor has a 32KB local scratchpad memory.
- `in[]`, `fa[]`, `ts[]`, and `lookup[]` fit in the local scratchpad memory close to the processor and can be read or written in a single cycle.
- For the baseline processor, `s[]` lives in the large (1MB) memory and requires 6 cycles to access.
- `lookup[]` and `s[]` are prepopulated with content before entering the while loop (not shown).
- Assume adds and multiplies take 1 ns when implemented in hardware accelerator, so fully pipelined accelerators also run at 1 GHz.

1. For sequential evaluation and assuming `FREQBYTES` is 256.

- (a) Worst-case cycles to compute one iteration of the outer while loop?
(show cycles per loop for partial credit consideration.)

A	WINDOW	1024
B	<code>FREQBYTES</code> × 10 100 MB/s bandwidth = 10 cycles/byte	2560
between	5	5
C	15 × <code>MAX_FREQS</code> 12 ops, 3 scratchpad memory accesses	3825
D, E	<code>MAX_FREQS</code> × (5 + WINDOW × 10) E 10: 6 for read from <code>s[]</code> + read and write <code>ts[]</code> + multiply, add	2,612,475
F	WINDOW	1024
Total		2,620,913

2.6 million cycles

- (b) Which outer loop is the bottleneck?

Circle One:

A	B	C	(D)	F
---	---	---	-----	---

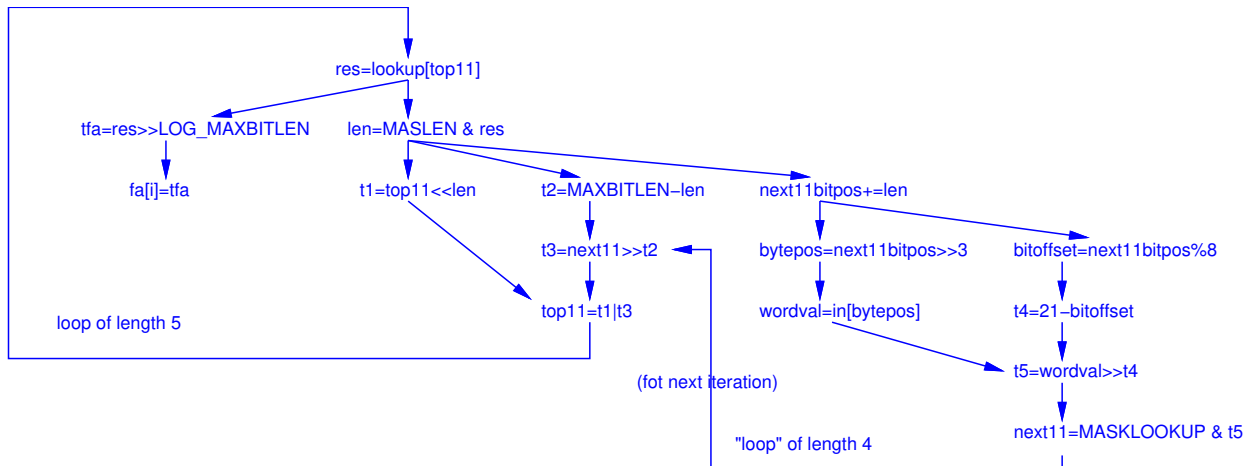
- (c) What is the Amdahl's Law maximum speedup for accelerating the identified loop?

$$\frac{A+B+C+D+F}{A+B+C+F} = \frac{2,620,913}{2,620,913 - 2,612,475} = 310$$

2. Loop C

- (a) How many memory operations does one instance of the loop perform?
 $3 - \text{lookup}[], \text{fa}[], \text{in}[]$
- (b) How many compute operations (of the set identified) does the loop perform?
 12
- (c) Assuming unlimited compute operators and memory ports, what is the minimum achievable Initiation Interval (II) for this loop?
 Draw dataflow graph and identify any data-dependent loops for full credit.

II=5



Note: Critical path is 7. The key loop is the one around `top11`, which is of length 5. We must also be able to update `next11`, and that is in a loop of length 4. Strictly, it's not a loop itself, but we do need to be able to compute the `next11` within one II, and this does fit.

3. Data Parallel: Classify Loops C, D, and E:

Loop	Data Parallel?	Associative Reduce?	Must be Sequential?
C			Yes
D		Yes	
E	Yes		

C: The dependent loop for `top11` identified in Problem2c forces sequentialization of the loop.

E: operations are independent for each j . Can perform the entire multiplication and add concurrently. This vectorizable.

D: If you think about unrolling E into a vector, then unrolling D as well, the only dependency is the add chain for each `freq` into `ts[j]`. The add is associative, so this is an associative reduce operation.

4. What is the latency bound for executing Loops C and D (from the beginning of C to the end of D)?

- assume memories of unbounded width (no bandwidth limits)
- respect latencies for memory access

Loop C: From Problems 2 and 3, we know this loop is sequentially dependent with an II of 5. So, it will take:
 $MAX_FREQS \times II = 255 \times 5 = 1275$ cycles.

Loop E: This is data parallel. Fully unrolled this takes 6 (read $s[][]$) + 1 = 7 cycles to get to the products.

Loop D: This is a reduce add across MAX_FREQS values to produce each $ts[j]$. That can be done in $\log_2(MAX_FREQS) = 8$ cycles. There's a final write into $ts[j]$ at the end.

Together, this gives us $1275 + 7 + 8 + 1 = 1291$ cycles or $1.3\mu s$.

We can do slightly better observing that we can overlap some (or most) of the additions in D-E with C. So, even if we sequentially perform the E vector adds, we can complete one per cycle and match pace with C. So, after finishing C, we only need to perform the 7 cycles for the E lookup and multiply, then a final add and store. So, we can perform this is $1275 + 7 + 1 + 1 = 1284$. To two significant figures, this is also $1.3\mu s$.

5. Data Streaming: How big (minimum size) does the buffer need to be between the identified loops in order to allow the loops to profitably execute concurrently. (Hint: Based on data dependencies, under what scenarios and granularity can the identified loops act as a producer-consumer pair in a pipeline.)

Explain size choices for partial credit consideration.

Loop Pair	Size (bytes)
B→C	1 or 4
C→D	4
D→F	4096

B→C: Each byte read can be passed directly to C, and C can perform a lookup. Technically, C may read a whole 32b word. However, depending on length, it may consume less than a byte on each iteration. If C is consuming less than a byte, it can use each byte as it shows up. If C is consuming a whole 32b word, then it will need to get a full word (4 bytes) to be able to perform each operation.

C→D: As each $fa[i]$ is produced, C can pass it to D, allowing D to perform one loop body on that $fa[i]$. fa is produced by C and consumed by D in order.

D→F: $ts[]$ is updated on every invocation of D. The final value of $ts[]$ is not known until the D completes the final iteration. As such, D cannot pass $ts[]$ to F until it completes its execution. Then the whole $ts[]$ ($WINDOW \times 4 = 4096$ bytes) can be given to F. F can write $ts[]$ out while D is operating on the next iteration of the outer while loop.

So, the whole B→C→D body can operate as a pipeline. B and C can operate on data in the same outer while iteration, passing data in bytes or words as they are produced, while F must operate on data from an earlier outer while iteration than B and C.

6. Consider trying to achieve a real-time rate of one window output per cycle (equivalently, the II of the outer while loop is WINDOW or 1024 cycles).

Assume you exploit data streaming between loops so they can run concurrently.

- (a) Given that Flash memory has a maximum throughput of 100 MB/s, what is the maximum possible value for `FREQBYTES`?

100MB/s throughput, means the fastest we can read each byte is once ever 10 cycles.

$$\text{FREQBYTES} \times 10 = 1024 \rightarrow \text{FREQBYTES}=102.$$

- (b) Based on your II identified in Problem 2c, what is the maximum value for `freqs` in order to meet this real-time throughput goal?

$$\text{freqs} \times \text{II} = 1024 \rightarrow \text{freqs} \times 5 = 1024 \rightarrow \text{freqs}=204.$$

- (c) What II do you need to achieve for Loop D to meet this real-time throughput goal?

The most direct argument is that this needs to match the rate of Loop C, so also has an II=5 requirement.

Alternately, we have the same equation, now with II_D as the variable.

$$\text{freqs} \times II_D = 1024 \rightarrow 204 \times II_D = 1024 \rightarrow II_D = 5.$$

7. Define the composition of a custom VLIW datapath for loop C that can achieve the identified II in Problem 2c.

For full credit, minimize area of your implementation.

Assume:

- Design includes at least one write port to a scratchpad memory containing fa[] and one read port to a scratchpad memory containing in[]
- Assume a crossbar interconnect between operator (and memory port) outputs and operator (and memory address, data) inputs.

- (a) How many operators of each type? Give both Resource Bound (RB) and number for which you can schedule.

Operator	Inputs	Outputs	Number	
			RB	Schedule
shifters	2	1	$\lceil \frac{5}{5} \rceil = 1$	2
ALU (includes , &, +, - , %-by-powers-of-2)	2	1	$\lceil \frac{7}{5} \rceil = 2$	2
scratchpad memory banks	2	1	1	1
ports to memory containing in[]	1	1	1	1
ports to memory containing fa[]	1	0	1	1
above error, should be	2	0		
branch units	1	0	1	1

- (b) How are the scratchpad memory banks used?

Hold lookup[] array.

- (c) Crossbar Inputs and Outputs for your design (final column, the one you can schedule)?

Inputs	13 (or 14 with correction)
Outputs	6

(d) Estimate the area for your design using the following costs.

- shifters: 1024
- ALU (includes |, &, +, -, %-by-powers-of-2): 32
- Scratchpad memory banks of depth d : $60(d + 6)$
- ports to memory containing in[]: 200
- ports to memory containing fa[]: 200
- branch unit: 100
- crossbar: $128 \times Inputs \times Outputs + 2400 \times Outputs$
(Each crossbar output includes a 4 word memory acting as a small register file for input to the associated operator or memory.)

$$2 \times 1024 + 2 \times 32 + 60(2048 + 6) + 200 + 200 + 100 + 128 \times 13 \times 6 + 2400 \times 6 = 150,236 \approx 150,000$$

(e) Provide a schedule:

Operator → Cycle	fa[] write	in[] read	Label with your selected operators						
			lookup[]	shift0	shift1	ALU0 ALU1 Branch			
0			res	t5					
1				tfa		len	next11		
2	fa[i]			t1		t2	next11bitpos		
3				t3	bytepos	bitoffset			
4		wordval				top11	t4	branch i < freqs	
5									
6									
7									
8									
9									
10									
11									
12									
13									
14									

Label cells with the variable assigned by the operation (or array entry written).

orange instructions software pipelined from previous iteration of loop

(Note extra schedules at end. May want to use as scratch while exploring schedules and put final here.)

8. Considering a custom hardware accelerator implementation where you are designing both the compute operators and the associated memory architecture, how would you use loop unrolling and array partitioning on Loop D to achieve the identified II in Problem 6c, while minimizing area?

Use the following area model and assume $s[]$, $ts[]$, and $fa[]$ are part of this loop module:

- n -bit counters: n
- 32b adder: 32
- 16×16 multiplier: 256
- Single-port, 32b-wide memory holding d words: $38(d + 6)$
- Double-port, 32b-wide memory holding d words: $60(d + 6)$

- (a) Unrolling for each loop (D, E)?

Loop	Unroll Factor
D	1
E	205

To meet the $II = 5$ goal, we must perform $\lceil \frac{1024}{5} \rceil = 205$ loop bodies of E on each cycle. So, we can unroll E 205 times and pipeline the computation.

- (b) For the unrolling, how many multipliers and adders?

Multipliers	205
Adders	205

Note: Since E is inside D, unrolling D D_{unroll} times and E E_{unroll} times, will result in $D_{unroll} \times E_{unroll}$ adders and multipliers.

- (c) Array partitioning for each array ($s[]$, $ts[]$, $fa[]$)?
(each memory block has either 1 or 2 ports)

Array	Array Partition	Ports (select one)		Words/partition
$s[]$	cyclic 205 dimension 1	(1)	2	1280
$ts[]$	cyclic 205	1	(2)	5
$fa[]$	1	1	(2)	256

Note that $s[]$ is only read. $ts[]$ must be both read and written on each iteration. $fa[]$ must be written by C and read by D.

Properly pipelined $fa[]$ could get away with one port; when C and D are run concurrently $fa[]$ could go away as a memory and just become a register between C and D.

- (d) Identify the component(s) that consumes most (>80%) of the area?
(you don't necessarily need to compute the area to fine precision, but you need to estimate where area is going well enough to answer the question above.)

Component	Calculate	Area
8-bit counter for E	8	11
3-bit counter for D	3	
Adder	32×205	6560
Multiplier	256×205	52480
$s[]$	$205 \times 38(1280 + 6)$	10017940
$ts[]$	$205 \times 60(5 + 6)$	135300
$fa[]$	$60(256 + 6)$	15720
total		10228011

98% of area is the single-ported memory for $s[]$.

Memory (for $s[]$) consumes >80% of the area.

This page left almost blank for pagination. You may use for answers and computations.

Extra schedule (in case you need it for trying schedules out, or if you need to put your answer here; be clear which schedule we should grade.)

Operator → Cycle	fa[] write	in[] read	Label with your selected operators															
0																		
1																		
2																		
3																		
4																		
5																		
6																		
7																		
8																		
9																		
10																		
11																		
12																		
13																		
14																		

Label cells with the variable assigned by the operation (or array entry written).

Extra schedule (in case you need it for trying schedules out, or if you need to put your answer here; be clear which schedule we should grade.)

Operator → Cycle	fa[] write	in[] read	Label with your selected operators																		
0																					
1																					
2																					
3																					
4																					
5																					
6																					
7																					
8																					
9																					
10																					
11																					
12																					
13																					
14																					

Label cells with the variable assigned by the operation (or array entry written).

Code of Academic Integrity

Since the University is an academic community, its fundamental purpose is the pursuit of knowledge. Essential to the success of this educational mission is a commitment to the principles of academic integrity. Every member of the University community is responsible for upholding the highest standards of honesty at all times. Students, as members of the community, are also responsible for adhering to the principles and spirit of the following Code of Academic Integrity.*

Academic Dishonesty Definitions

Activities that have the effect or intention of interfering with education, pursuit of knowledge, or fair evaluation of a students performance are prohibited. Examples of such activities include but are not limited to the following definitions:

A. Cheating Using or attempting to use unauthorized assistance, material, or study aids in examinations or other academic work or preventing, or attempting to prevent, another from using authorized assistance, material, or study aids. Example: using a cheat sheet in a quiz or exam, altering a graded exam and resubmitting it for a better grade, etc.

B. Plagiarism Using the ideas, data, or language of another without specific or proper acknowledgment. Example: copying another persons paper, article, or computer work and submitting it for an assignment, cloning someone elses ideas without attribution, failing to use quotation marks where appropriate, etc.

C. Fabrication Submitting contrived or altered information in any academic exercise. Example: making up data for an experiment, fudging data, citing nonexistent articles, contriving sources, etc.

D. Multiple Submissions Multiple submissions: submitting, without prior permission, any work submitted to fulfill another academic requirement.

E. Misrepresentation of academic records Misrepresentation of academic records: misrepresenting or tampering with or attempting to tamper with any portion of a students transcripts or academic record, either before or after coming to the University of Pennsylvania. Example: forging a change of grade slip, tampering with computer records, falsifying academic information on ones resume, etc.

F. Facilitating Academic Dishonesty Knowingly helping or attempting to help another violate any provision of the Code. Example: working together on a take-home exam, etc.

G. Unfair Advantage Attempting to gain unauthorized advantage over fellow students in an academic exercise. Example: gaining or providing unauthorized access to examination materials, obstructing or interfering with another students efforts in an academic exercise, lying about a need for an extension for an exam or paper, continuing to write even when time is up during an exam, destroying or keeping library materials for ones own use., etc.

* If a student is unsure whether his action(s) constitute a violation of the Code of Academic Integrity, then it is that students responsibility to consult with the instructor to clarify any ambiguities.