

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

ESE532, Fall 2017

Final

Thursday, December 21

- Exam ends at 5:00PM; begin as instructed (target 3:00PM)
- 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.

<b>Name:</b>
--------------

Problem 1			Problem 2			Problem 3		Problem 4					Total
a	b	c	a	b	c	a	b	a	b	c	d	e	
10	10	5	10	5	10	15	10	5	5	5	10	5	100

1. Consider the following set of real-time tasks for a vision-based control system:

Task	Frequency	Operations per Invocation
Plan	once per second	$1.28 \times 10^{11}$
Control	100 times per second	$10^6$
Mapping	50 times per second	$4 \times 10^7$
Collide	100 times per second	$10^7$

Consider the following computational elements:

Compute Element	Performance	Area
High Performance Processor (HPP)	1 op per cycle at 2 GHz	$3 \times 10^4$
VLIW	8 ops per cycles at 500 MHz	$3 \times 10^3$
6-LUT (with interconnect)	200 MHz	1

Consider the following spatial datapaths for each of the tasks:

Task	6-LUTs	Effective operations per cycle
Plan	1,300	40
Control	32,000	500
Mapping	33,000	1024
Collide	3,200	100

Assume:

- perfectly divide work among processors (if need more than one),
- perfectly schedule VLIW
- perfectly divide task among instances of spatial datapaths shown above
- may share processor or VLIW among multiple tasks

- (a) Complete table with the number of homogeneous resources required to meet real-time task goals and the associated area: (May share processors, VLIW, but must pay for an integer number of each.)

Task	HPP	VLIW	6-LUTs
Plan			
Control			
Mapping			
Collide			
Total (integer)			
Total Area			

- (b) Heterogeneous solution that meets real-time goals and minimizes area. For each task identify the number and type of compute units.  
 (May share processors, VLIW, but must pay for an integer number of each. Show number of datapaths for 6-LUT case, and calculate total area of those datapaths for Area column, when appropriate.)

Task	Type	Number	Area
Plan			
Control			
Mapping			
Collide			
Total Area			

- (c) Show schedule for any HPP or VLIW processors that run multiple tasks at the coarsest time-slice for the task set.

## 2. Parallelism in Plan Task

Working with the following algorithm for Plan and assuming this takes the  $1.28 \times 10^{11}$  operations stated in Problem 1:

```
#define DIM 1024
typedef struct cell {
    uint32_t cost;
    uint32_t path;
    uint16_t edge_cost[4];
} cell_struct;
cell_struct from[DIM][DIM];
cell_struct to[DIM][DIM];
int xoff, yoff;

for (int iter=0;iter<3*DIM;iter++) {
    for (int y=0;y<DIM;y++)
        for (int x=0;x<DIM;x++) {
            uint32_t min_cost=from[y][x].cost;
            uint32_t min_path=from[y][x].path;
            for (int j=0;j<4;j++) {
                if (j==0) {xoff=0; yoff=-1;}
                else if (j==1) {xoff=-1; yoff=0;}
                else if (j==2) {xoff=1; yoff=0;}
                else if (j==3) {xoff=0; yoff=1;}
                uint32_t ncost=from[y+yoff][x+xoff].cost;
                uint32_t ecost=from[y][x].edge_cost[j];
                uint32_t icost=ncost+ecost
                if (icost < min_cost){
                    min_path=ENCODE(yoff,xoff); // count 3 ops (depth 1 op)
                    min_cost=icost;
                } // icost < min_cost
            } // for j
            to[y][x].cost=min_cost;
            to[y][x].path=min_path;
        } // for x
    cell_struct **tmp=from;
    from=to;
    to=tmp;
} // for iter
```

- (a) Describe how you can divide the computational task among processors or hardware datapaths to achieve various levels of parallelism. Write code snippets or edits as necessary to make your division clear.

- (b) Building on your answer above, how do you provide the specific parallelism needed in Problem 1? (Number is the value you calculated in Problem 1)

Compute Element	Number	How
HPP		
VLIW		
6-LUT datapaths		

- (c) What is the Latency Bound (critical path) in operations for the entire task assuming the data needed lives in registers (i.e., no latency coming or going from memory)?

## 3. Hardware datapath for Plan

Working with the same algorithm for Plan:

```

#define DIM 1024
typedef struct cell {
    uint32_t cost;
    uint32_t path;
    uint16_t edge_cost[4];
} cell_struct;
cell_struct from[DIM][DIM];
cell_struct to[DIM][DIM];
int xoff, yoff;

for (int iter=0;iter<3*DIM;iter++) {
    for (int y=0;y<DIM;y++)
        for (int x=0;x<DIM;x++) {
            uint32_t min_cost=from[y][x].cost;
            uint32_t min_path=from[y][x].path;
            for (int j=0;j<4;j++) {
                if (j==0) {xoff=0; yoff=-1;}
                else if (j==1) {xoff=-1; yoff=0;}
                else if (j==2) {xoff=1; yoff=0;}
                else if (j==3) {xoff=0; yoff=1;}
                uint32_t ncost=from[y+yoff][x+xoff].cost;
                uint32_t ecost=from[y][x].edge_cost[j];
                uint32_t icost=ncost+ecost
                if (icost < min_cost){
                    min_path=ENCODE(yoff,xoff); // count 3 ops (depth 1 op)
                    min_cost=icost;
                } // icost < min_cost
            } // for j
            to[y][x].cost=min_cost;
            to[y][x].path=min_path;
        } // for x
    cell_struct **tmp=from;
    from=to;
    to=tmp;
} // for iter

```

- (a) Unroll and pipeline the inner (j) loop.  
Draw the resulting pipelined datapath.

- (b) Show how you manage data storage and movement of **from** data in order to deliver **from** data to the datapath above. Assume data is streamed from the large memory. Draw registers and memories and show inputs to the Part a datapath



(This page intentionally left mostly blank for pagination.  
You may use for extra answer or work space.)

## 4. Data Management

Continuing to work with the algorithm for Plan:

```

#define DIM 1024
typedef struct cell {
    uint32_t cost;
    uint32_t path;
    uint16_t edge_cost[4];
} cell_struct;
cell_struct from[DIM][DIM];
cell_struct to[DIM][DIM];
int xoff, yoff;

for (int iter=0;iter<3*DIM;iter++) {
    for (int y=0;y<DIM;y++)
        for (int x=0;x<DIM;x++) {
            uint32_t min_cost=from[y][x].cost;
            uint32_t min_path=from[y][x].path;
            for (int j=0;j<4;j++) {
                if (j==0) {xoff=0; yoff=-1;}
                else if (j==1) {xoff=-1; yoff=0;}
                else if (j==2) {xoff=1; yoff=0;}
                else if (j==3) {xoff=0; yoff=1;}
                uint32_t ncost=from[y+yoff][x+xoff].cost;
                uint32_t ecost=from[y][x].edge_cost[j];
                uint32_t icost=ncost+ecost
                if (icost < min_cost){
                    min_path=ENCODE(yoff,xoff); // count 3 ops (depth 1 op)
                    min_cost=icost;
                } // icost < min_cost
            } // for j
            to[y][x].cost=min_cost;
            to[y][x].path=min_path;
        } // for x
    cell_struct **tmp=from;
    from=to;
    to=tmp;
} // for iter

```

- Processor has a local memory that holds 64KB with single cycle access for 16b, 32b, or 64b data.
  - `uint32_t` is stored in 4 Bytes; `uint16_t` is stored in 2 Bytes.
  - `from` and `to` live in the large memory.
  - Large memory has a read latency of 50 ns (100 cycles on HPP) for a single 16b, 32b, or 64b word or 100 ns (200 cycles) for a 2048b transfer into the local memory
    - For the processor, assume you have a routine
 

```
void read2048(int *local_buf, int *large_mem_buf)
```

 that will initiate a 2048b read from the large memory into the processor local memory at the specified addresses.
  - The processor can issue a write to the large memory in a single cycle and continue without waiting.
    - We should be concerned with when a read may need something previously written to memory and the impact of writes on memory bandwidth. To keep this problem simple, you may ignore these effects.
  - Assume scalar (non-array) variables can live in registers.
  - You may ignore loop and conditional overheads in processor runtime estimates; additions and comparisons are single cycle operations.
- (a) What time (in ns) is required for the task as written when running on the HPP and all references to each component of `from` go to the large memory?
- (b) What fraction of time is spent performing read operations from the large memory for `from`?
- (c) What is the Amdahl's Law speedup if you only accelerate these memory reads to `from`?

- (d) How do you need to modify the design to minimize the time spent reading from the large memory?  
Describe how the code needs to be modified. Show revised code as necessary for clarity.

- (e) What is the time (in ns) required on a single HPP processor after this modification?

## 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.