

ETH 263-2210-00L COMPUTER ARCHITECTURE, FALL 2019

HW 5: GPU PROGRAMMING, EMERGING MEMORY TECHNOLOGIES, PREFETCHING
ASYMMETRIC MULTIPROCESSORS, CONSISTENCY, AND COHERENCE

Instructor: Prof. Onur Mutlu

TAs: Mohammed Alser, Rahul Bera, Geraldo Francisco De Oliveira Junior, Can Firtina,
Juan Gomez Luna, Jawad Haj-Yahya, Hasan Hassan, Konstantinos Kanellopoulos, Jeremie Kim,
Nika Mansouri Ghiasi, Lois Orosa Nogueira, Jisung Park, Minesh Hamenbhai Patel, Abdullah Giray Yaglikci

Given: Friday, Dec 6, 2019

Due: **Thursday, Dec 19, 2019**

- **Handin - Critical Paper Reviews (1).** You need to submit your reviews to <https://safari.ethz.ch/review/architecture19/>. Please, check your inbox, you should have received an email with the password you should use to login. If you didn't receive any email, contact comparch@lists.inf.ethz.ch. In the first page after login, you should click in "Architecture - Fall 2019 Home", and then go to "any submitted paper" to see the list of papers.
- **Handin - Questions (2-7).** You should upload your answers to the Moodle Platform (<https://moodle-app2.let.ethz.ch/mod/assign/view.php?id=391622>) as a single PDF file.

1. Critical Paper Reviews [300 points]

Please read the guidelines for reviewing papers and check the sample reviews. You may access them by *simply clicking on the QR codes below or scanning them*. We will give out extra credit that is worth 0.5% of your total grade for each good review.



Guidelines



Sample reviews

Write an approximately one-page critical review for each of the following papers. A review with bullet point style is more appreciated. Try not to use very long sentences and paragraphs. Keep your writing and sentences simple. Make your points bullet by bullet, as much as possible.

- Suleman et al., "Accelerating Critical Section Execution with Asymmetric Multi-Core Architectures" in Proceedings of the 14th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2009. https://people.inf.ethz.ch/omutlu/pub/acs_asplos09.pdf
- Srinath et al., "Feedback Directed Prefetching: Improving the Performance and Bandwidth-Efficiency of Hardware Prefetchers" in Proceedings of the 13th International Symposium on High-Performance Computer Architecture (HPCA), 2006. https://people.inf.ethz.ch/omutlu/pub/srinath_hPCA07.pdf
- Vijaykumar et al., "A Case for Core-Assisted Bottleneck Acceleration in GPUs: Enabling Flexible Data Compression with Assist Warps" in Proceedings of the 42nd International Symposium on Computer Architecture (ISCA), 2015. https://people.inf.ethz.ch/omutlu/pub/caba-gpu-assist-warps_isca15.pdf

2. GPU Programming and Performance Analysis [150 points]

The following two program segments are executed on a GPU with C compute units. In each compute unit, one or more thread-blocks can run. Each thread-block is composed of threads that are grouped into warps of W threads.

In both programs, 2D thread-blocks are used. Each thread-block is identified by its block indices (bx , by), and each thread is identified by its thread indices (tx , ty). The size of a thread-block is $bdx * bdy$. Consider that a thread-block is decomposed into warps in a way that threads with consecutive tx and equal ty belong to the same warp. More specifically, the warp number for a thread (tx , ty) is $\frac{ty * bdx + tx}{warp-size}$.

The entire input size is $rows * cols$ integers. The size of an integer element is 4 bytes. The input is divided into tiles that are assigned to the thread-blocks.

`local_data` is an array in *local memory*, a fast on-chip memory that is used as a software-programmable cache. The amount of local memory per compute unit is S bytes. The threads of a thread-block can load data from *global memory* (i.e., the GPU off-chip memory) into local memory. The size of a global memory transaction is equal to the warp size times 4 bytes.

Program A:

```
__gpu_kernel_a(int* data, int rows, int cols){

    int* local_data[bdx * bdy];

    const int g_row = by * bdy + ty;
    const int g_col = bx * bdx + tx;
    const int l_row = ty;
    const int l_col = tx;
    const int pos    = g_row * cols + g_col;

    local_data[l_row * bdx + l_col] = data[pos];

    // Compute using local_data
}
```

Program B:

```
__gpu_kernel_b(int* data, int rows, int cols){

    int* local_data[bdx * bdy];

    const int g_row = bx * bdx + tx;
    const int g_col = by * bdy + ty;
    const int l_row = tx;
    const int l_col = ty;
    const int pos    = g_row * cols + g_col;

    local_data[l_row * bdy + l_col] = data[pos];

    // Compute using local_data
}
```

Please answer the questions on the next page.

(a) What is the maximum number of thread-blocks that run in *each* compute unit for programs A and B?

A large, empty rectangular box with a thin black border, intended for the student to write their answer to question (a).

(b) Assuming that the GPU does *not* have caches, which program will execute faster? Why?

A large, empty rectangular box with a thin black border, intended for the student to write their answer to question (b).

(c) Assume that the GPU has a single level of cache shared by all compute units. What will be the effect of this cache on the execution time of programs A and B?

A large, empty rectangular box with a thin black border, intended for the student to write their answer to question (c).

- (d) Assume that the access latency to the shared cache in part (c) is negligible. What should be the minimum size of the shared cache to guarantee that programs A and B have the same (or very similar) performance? (NOTE: The solution is independent of the warp scheduling policy).



- (e) Now assume that *only one* thread-block is executed in each compute unit. Each thread-block in program A needs always T ms to complete its work, because the computation is very regular. What will be the total execution time of program A?



3. Emerging Memory Technologies [150 points]

Computer scientists at ETH developed a new memory technology, ETH-RAM, which is non-volatile. The access latency of ETH-RAM is close to that of DRAM while it provides higher density compared to the latest DRAM technologies. ETH-RAM has one shortcoming, however: it has limited endurance, i.e., a memory cell stops functioning after 10^6 writes are performed to the cell (known as cell wear-out).

A bright ETH student has built a computer system using 1 GB of ETH-RAM as main memory. ETH-RAM exploits a perfect wear-leveling mechanism, i.e., a mechanism that equally distributes the writes over all of the cells of the main memory.

- (a) This student is worried about the lifetime of the computer system she has built. She executes a test program that runs special instructions to bypass the cache hierarchy and repeatedly writes data into different words until **all** the ETH-RAM cells are worn-out (stop functioning) and the system becomes useless. The student's measurements show that ETH-RAM stops functioning (i.e., all its cells are worn-out) in one year (365 days). Assume the following:
- The processor is in-order and there is no memory-level parallelism.
 - It takes 5 ns to send a memory request from the processor to the memory controller and it takes 28 ns to send the request from the memory controller to ETH-RAM.
 - ETH-RAM is word-addressable. Thus, each write request writes 4 bytes to memory.

What is the write latency of ETH-RAM? Show your work.

- (b) ETH-RAM works in the multi-level cell (MLC) mode in which each memory cell stores 2 bits. The student decides to improve the lifetime of ETH-RAM cells by using the single-level cell (SLC) mode. When ETH-RAM is used in SLC mode, the lifetime of each cell improves by a factor of 10 and the write latency decreases by 70%. What is the lifetime of the system using the SLC mode, if we repeat the experiment in part (a), with everything else remaining the same in the system? Show your work.

4. Prefetching [150 points]

A processor is observed to have the following access pattern to cache blocks. Note that the addresses are **cache block addresses**, not byte addresses. This pattern is repeated for a large number of iterations.

Access Pattern P: A , $A + 3$, $A + 6$, A , $A + 5$

Each cache block is 8KB. The hardware has a fully associative cache with LRU replacement policy and a total size of 24KB.

None of the prefetchers mentioned in this problem employ confidence bits, but they all start out with empty tables at the beginning of the access stream shown above. Unless otherwise stated, assume that 1) each access is separated long enough in time such that all prefetches issued can complete before the next access happens, and 2) the prefetchers have large enough resources to detect and store access patterns.

- (a) You have a stream prefetcher (i.e., a next- N -block prefetcher), but you don't know the prefetch degree (N) of it. However, you have a magical tool that displays the coverage and accuracy of the prefetcher. When you run a large number of repetitions of access pattern P , you get 40% coverage and 10% accuracy. What is the degree of this prefetcher (how many next blocks does it prefetch)?

- (b) You didn't like the performance of the stream prefetcher, so you switched to a PC-based stride prefetcher that issues prefetch requests based on the stride detected for each memory instruction. Assume all memory accesses are incurred by the *same* load instruction (i.e., the same PC value) and the initial stride value for the prefetcher is set to 0.

Circle which of the cache block addresses are prefetched by this prefetcher:

A , $A + 3$, $A + 6$, A , $A + 5$
 A , $A + 3$, $A + 6$, A , $A + 5$
 A , $A + 3$, $A + 6$, A , $A + 5$
 A , $A + 3$, $A + 6$, A , $A + 5$

Explain:

- (c) Stride prefetcher couldn't satisfy you either. You changed to a Markov prefetcher with a correlation table of 12 entries (assume each entry can store a single address to prefetch, and remembers the most recent correlation). When all the entries are filled, the prefetcher replaces the entry that is least-recently accessed.

Circle which of the cache block addresses are prefetched by this prefetcher:

A , $A + 3$, $A + 6$, A , $A + 5$
 A , $A + 3$, $A + 6$, A , $A + 5$
 A , $A + 3$, $A + 6$, A , $A + 5$
 A , $A + 3$, $A + 6$, A , $A + 5$

Explain:

- (d) Just in terms of coverage, after how many repetitions of access pattern P does the Markov prefetcher from part (c) start to outperform the stream prefetcher from part (a), if it can at all? Show your work.

- (e) You think having a correlation table of 12 entries makes the hardware too costly, and want to reduce the number of correlation table entries for the Markov prefetcher. What is the minimum number of entries that gives the same prefetcher performance as 12 entries? Similar to the last part, assume each entry can store a single next address to prefetch, and remembers the most recent correlation. Show your work.

- (f) Your friend is running the same program on a different machine that has a Markov prefetcher with 2 entries. The same assumptions from part (e) apply.

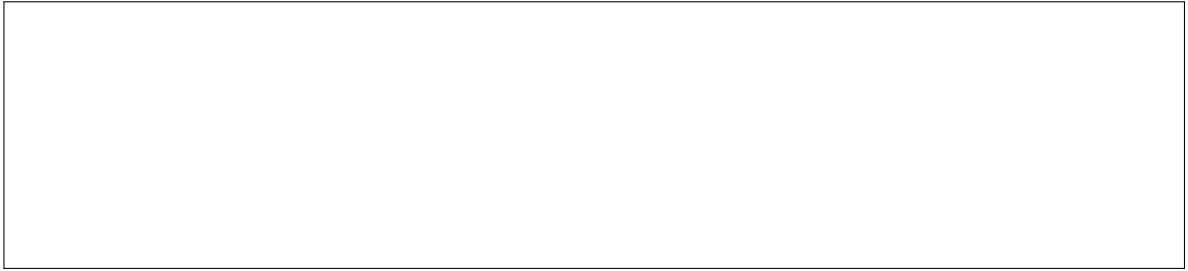
Circle which of the cache block addresses are prefetched by the prefetcher:

A, A + 3, A + 6, A, A + 5
A, A + 3, A + 6, A, A + 5
A, A + 3, A + 6, A, A + 5
A, A + 3, A + 6, A, A + 5

Explain:

- (g) As an avid computer architect, you decide to update the processor by increasing the cache size to 32KB with the same cache block size. Assume you will be only running a program with the same access pattern P for a large number of iterations (i.e., one trillion), describe a prefetcher that provides

smaller memory bandwidth consumption than the baseline without a prefetcher.
Explain:

A large empty rectangular box with a thin black border, intended for the user to provide an explanation for the statement above.

5. Asymmetric Multicore [150 points]

A microprocessor manufacturer asks you to design an asymmetric multicore processor for modern workloads. You should optimize it assuming a workload with 80% of its work in the parallel portion. Your design contains one large core and several small cores, which share the same die. Assume the total die area is 32 units.

- *Large core:* For a large core that is n times faster than a single small core, you will need n^3 units of die area (n is a positive integer). The dynamic power of this core is $6 \times n$ Watts and the static power is n Watts.
- *Small cores:* You will fit as many small cores as possible, after placing the large core. A small core occupies 1 unit of die area. Its dynamic power is 1 Watt and its static power is 0.5 Watts.

The parallel portion executes *only* on the small cores, while the serial portion executes *only* on the large core.

Please answer the following questions. Show your work. Express your equations and solve them. You can approximate some computations, and get partial or full credit.

- (a) What configuration (i.e., number of small cores and size of the large core) results in the best performance?

- (b) The energy consumption should also be a metric of reference in your design. Compute the energy consumption for the best configuration in part (a).



- (c) For the best configuration obtained in part (a), you are considering to use the large core to collaborate with the small cores on the execution of the parallel portion.

- (i) What is the overall performance improvement, compared to the performance obtained in part (a), if the large core collaborates on the parallel portion?



- (ii) What is the overall energy change, compared to the energy obtained in part (b), if the large core collaborates on the parallel portion?



- (iii) Discuss whether it is worth using the large core to collaborate with the small cores on the execution of the parallel portion.



- (d) Now assume that the serial portion can be optimized, i.e., the serial portion becomes smaller. This gives you the possibility of reducing the size of the large core, and still improving performance. For a large core with an area of $(n - 1)^3$, where n is the value obtained in part (a), what should be the fraction of serial portion that would lead to better performance than in part (a)?

- (e) Your design is so successful for desktop processors that the company wants to produce a similar design for mobile devices. The power budget becomes a constraint. For a maximum of total power of 20W, how much would you need to reduce the dynamic power consumption of the large core, if at all, for the best configuration obtained in part (a)? Assume again that the parallel fraction is 80% of the workload. (Hint: Express the dynamic power of the large core as $D \times n$ Watts, where D is a constant).

6. Memory Consistency [150 points]

A programmer writes the following two C code segments. She wants to run them concurrently on a multicore processor, called SC, using two different threads, each of which will run on a different core. The processor implements *sequential consistency*, as we discussed in the lecture.

Thread T0		Thread T1	
Instr. T0.0	<code>a = X[0];</code>	Instr. T1.0	<code>Y[0] = 1;</code>
Instr. T0.1	<code>b = a + Y[0];</code>	Instr. T1.1	<code>*flag = 1;</code>
Instr. T0.2	<code>while(*flag == 0);</code>	Instr. T1.2	<code>X[1] *= 2;</code>
Instr. T0.3	<code>Y[0] += 1;</code>	Instr. T1.3	<code>a = 0;</code>

`X`, `Y`, and `flag` have been allocated in main memory, while `a` and `b` are contained in processor registers. A read or write to any of these variables generates a single memory request. The initial values of all memory locations and variables are 0. Assume each line of the C code segment of a thread is a *single* instruction.

- (a) What is the final value of `Y[0]` in the SC processor, after both threads finish execution? Explain your answer.

- (b) What is the final value of `b` in the SC processor, after both threads finish execution? Explain your answer.

With the aim of achieving higher performance, the programmer tests her code on a new multicore processor, called RC, that implements *weak consistency*. As discussed in the lecture, the weak consistency model has no need to guarantee a strict order of memory operations. For this question, consider a very weak model where there is *no* guarantee on the ordering of instructions as seen by different cores.

- (c) What is the final value of $Y[0]$ in the RC processor, after both threads finish execution? Explain your answer.

After several months spent debugging her code, the programmer learns that the new processor includes a `memory_fence()` instruction in its ISA. The semantics of `memory_fence()` is as follows for a given thread that executes it:

1. Wait (stall the processor) until *all* preceding memory operations from the thread complete in the memory system and become visible to other cores.
2. Ensure *no* memory operation from any later instruction in the thread gets executed before the `memory_fence()` is retired.

(d) What *minimal* changes should the programmer make to the program above to ensure that the final value of `Y[0]` on RC is the same as that in part (a) on SC? Explain your answer.

7. Cache Coherence [150 points]

We have a system with 4 byte-addressable processors. Each processor has a private 256-byte, direct-mapped, write-back L1 cache with a block size of 64 bytes. Coherence is maintained using the Illinois Protocol (MESI), which sends an invalidation to other processors on writes, and the other processors invalidate the block in their caches if *the block is present* (NOTE: On a write hit in one cache, a cache block in Shared state becomes Modified in that cache).

Accessible memory addresses range from 0x50000000 – 0x5FFFFFFF. Assume that the offset within a cache block is 0 for all memory requests. We use a snoopy protocol with a shared bus.

Cosmic rays strike the MESI state storage in your coherence modules, causing the state of a *single* cache line to instantaneously change to another state. This change causes an inconsistent state in the system. We show below the initial tag store state of the four caches, *after* the inconsistent state is induced.

Initial State

Cache 0		
	Tag	MESI state
Set 0	0x5FFFFFFF	M
Set 1	0x5FFFFFFF	E
Set 2	0x5FFFFFFF	S
Set 3	0x5FFFFFFF	I

Cache 1		
	Tag	MESI state
Set 0	0x522222	I
Set 1	0x510000	S
Set 2	0x5FFFFFFF	S
Set 3	0x533333	S

Cache 2		
	Tag	MESI state
Set 0	0x5F111F	M
Set 1	0x511100	E
Set 2	0x5FFFFFFF	S
Set 3	0x533333	S

Cache 3		
	Tag	MESI state
Set 0	0x5FF000	E
Set 1	0x511100	S
Set 2	0x5FFFF0	I
Set 3	0x533333	I

- (a) What is the inconsistency in the above initial state? Explain with reasoning.

(b) Consider that, after the initial state, there are several paths that the program can follow that access different memory instructions. In b.1-b.4, we will examine whether the followed path can potentially lead to incorrect execution, i.e., an incorrect result.

b.1) Could the following path potentially lead to incorrect execution? Explain.

order	Processor 0	Processor 1	Processor 2	Processor 3
1 st			ld 0x51110040	
2 nd	st 0x5FFFFFF40			
3 rd				st 0x51110040
4 th		ld 0x5FFFFFF80		
5 th		ld 0x51110040		
6 th		ld 0x5FFFFFF40		

b.2) Could the following path potentially lead to incorrect execution? Explain.

order	Processor 0	Processor 1	Processor 2	Processor 3
1 st				ld 0x51110040
2 nd	ld 0x5FFFFFF00			
3 rd			ld 0x51234540	
4 th	st 0x5FFFFFF40			
5 th				ld 0x51234540
6 th	ld 0x5FFFFFF00			

After some time executing a particular path (which could be a path *different* from the paths in parts b.1-b.4) and with no further state changes caused by cosmic rays, we find that the final state of the caches is as follows.

Final State

Cache 0		
	Tag	MESI state
Set 0	0x5FFFFFFF	M
Set 1	0x5FFFFFFF	E
Set 2	0x5FFFFFFF	S
Set 3	0x5FFFFFFF	E

Cache 1		
	Tag	MESI state
Set 0	0x5FF000	I
Set 1	0x510000	S
Set 2	0x5FFFFFFF	S
Set 3	0x533333	I

Cache 2		
	Tag	MESI state
Set 0	0x5F111F	M
Set 1	0x511100	E
Set 2	0x5FFFFFFF	S
Set 3	0x533333	I

Cache 3		
	Tag	MESI state
Set 0	0x5FF000	M
Set 1	0x511100	S
Set 2	0x5FFFF0	I
Set 3	0x533333	I

- (c) What is the *minimum* set of memory instructions that leads the system from the initial state to the final state? Indicate the set of instructions in order, and clearly specify the access type (ld/st), the address of each memory request, and the processor from which the request is generated.