

Homework #4

Q1.

(PLEASE NOTE: Matheus said in his office hours (YouTube video on Wednesday 3/11/15) that filling out the table is essentially showing our work. Based on that assumption, I didn't find it necessary to type all my thought process for this question (Figuring out how to type up our thought process in English is just a challenging task on this quesiton) PLEASE don't deduct points for not showing my thought process.)

	Core	Instr	Addr	L1 Access	Current L1 State	L1-L2 Bus	Current L2 State	L2 Access	L2 State	L1 State
1	0	Ld	0x100	Miss	I	RRd	I	Miss	E	E
2	0	Ld	0x150	Miss	I	RRd	I	Miss	E	E
3	1	St	0x100	Miss	I	RRx	E	Hit	M	M
4	1	St	0x104	Hit	M	- *	M	- *	M	M
5	1	Ld	0x150	Miss	I	RRd	E	Hit	E	S
6	2	Ld	0x154	Miss	I	RRd	I	Miss	S	S
7	3	St	0x100	Miss	I	RRx/Bwb	I	Miss	M	M
8	3	St	0x150	Miss	I	RRx	S	Hit	M	M
9	1	Ld	0xA00100	Miss	I	RRd	I	Miss	E	E
10	1	Ld	0xB00100	Miss	I	RRd	I	Miss	E	E
11	1	Ld	0x100	Miss	I	RRd	I	Miss	S	E

*NOTE: Since we have a hit at L1, no requests will be send out to L1-L2 Bus.

Q2A.

This is for one cacheline

Addr :	a[0][0]	miss
Addr+0x4:	a[0][1]	hit
Addr+0x8:	a[0][2]	hit
.		hit
.		

This is for the second cacheline

		miss
		hit
		hit
Addr+...	a[0][7]	hit

total of 2 cachelines => 2 misses for each

Sine first access at the beginning of each cache considers being a cold miss, so we would have a miss per each cacheline. Therefore we have 2 misses for every 2 cachelines. We would have the following situation:
(assumption: I'm computing number of misses at L1)

As we could see it does matter how threads are mapped to each processor.

Q2B.

Hit rate will decrease, because based on the mapping given there is less spatial locality. We have the situation below:

Addr+0x0: a[0][0]

Addr+0x10: a[0][1]

Addr+0x20: a[0][2]

Addr+0x30: a[0][3]

Addr+0x40: a[0][4]

Addr+0x50: a[0][5]

Addr+0x60: a[0][6]

Addr+0x70: a[0][7]

.

.

.

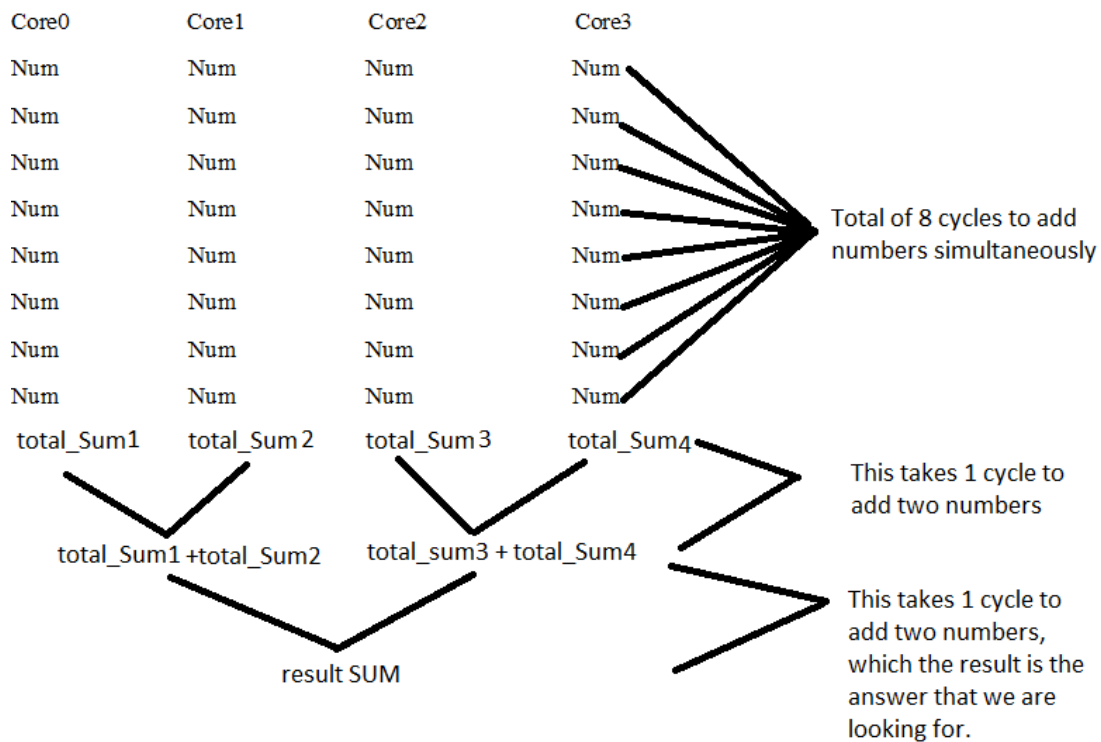
There is going to be gaps between addresses, which it requires more time to be found (or to be accessed).

Q2.C

In single thread program the number of cache hits don't change. But, if we are dealing with multi-threaded program, it really depends on how we arrange the threads (like the case we had in previous part, the number of misses was depended on how we mapped the threads to cores) and they overlap, so the order of accesses to cache hierarchy and memory changes could be arranged differently, which we could end up getting different hit rate (or miss rate).

Q2.D

In order to parallelize the operation of sum, we could do the operation of sum using 4 cores, which they add 8 number each. We have the situation below:



Therefore we will have the final sum by $8(\text{cycles}) + 1(\text{cycle}) + 1(\text{cycle}) = 10$ cycles in total

In order to convert the single threaded code to multi-threaded code, we do the following:

```
int reduction ( int **c , int tid) {
```

```
    int sum = 0;
```

```
    for (int j = 0; j < 8; j++){
```

```
        sum += c[tid][j];
```

```
    }
```

```
    return sum;
```

```
}
```

```
void main(){
```

```
    .
```

```
    .
```

```
    .
```

```
// we will have total of 8 add (assumption: we only care about add instruction) this logic sort of follows the
```

```
// diagram shown above. We have 4 threads in which each does 8 additions.
```

```
reduction_sum[0] = reduction(c , 0);
```

```
reduction_sum[1] = reduction(c , 1);
```

```
reduction_sum[2] = reduction(c , 2);
```

```
reduction_sum[3] = reduction(c , 3);
```

```
// This is where we add the 4 result computed by each thread.

reduction_Sum0 = reduction_sum[0] + reduction_sum[1]; // This is adding the value of sum found by core0 and core1
reduction_Sum1 = reduction_sum[2] + reduction_sum[3]; // This is adding the value of sum found by core2 and core3


// At the end we add up the two sum values computed by each two cores:

result_Sum = reduction_Sum0 + reduction_Sum1;          // computing the value of total sum.

.
.
.

return 0;

}
```