

Cache Awareness

Course Level:

CS1/CS2

PDC Concepts Covered:

PDC Concept	Bloom Level
Locality	C
False Sharing	C

Programming Knowledge Prerequisites:

- Know how to compile C/C++
- Be able to understand loops and functions

Tools Required:

For C/C++: A C++ compiler that is OpenMP capable (e.g. the Gnu gcc C++ compiler)

Activity:

Consider the code below. It creates an array of size $SIZE * SIZE$ and then loops over the array setting each element to zero. The time it takes to loop over the array is measured and printed out. The loop is implemented with two for loops, each iterating from 0 to $SIZE$. This lets us think of the array as a matrix with $SIZE$ rows and $SIZE$ columns. Then i represents the row index and j represents the column index. The expression $i * SIZE + j$ calculates the actual index in the array. As we can see from the loops, the matrix is zeroed row by row.

Compile and run the following code in C++:

```
1 #include <iostream>
2 #include "sys/time.h"
3
4 using namespace std;
5
6 int main()
7 {
8     const int SIZE = 20000;
9     int* array = new int[SIZE*SIZE];
10    struct timeval tv1, tv2;
11
12    gettimeofday(&tv1, NULL); //get the current time
13
14    for(int i = 0; i < SIZE; i++)
15        for(int j = 0; j < SIZE; j++)
```

```

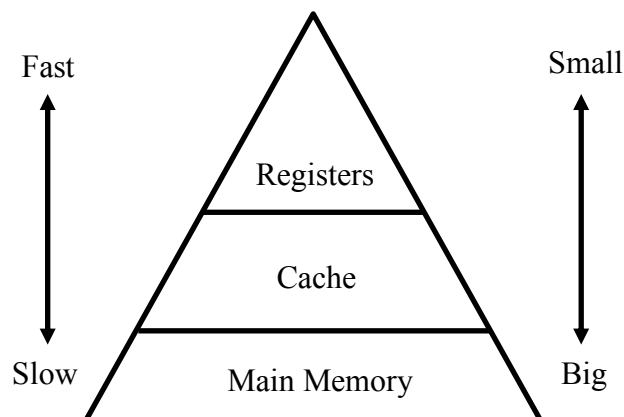
16         array[i*SIZE+j]=0;
17
18         gettimeofday(&tv2,NULL);
19
20         //this prints out the time in seconds between
21         //the two calls to gettimeofday
22         cout << "Time: " << (double)(tv2.tv_usec - tv1.tv_usec) / 1000000 +
23         (double)(tv2.tv_sec - tv1.tv_sec) << endl;
24
25         delete [] array;
26         return 0;
27     }

```

What if we want to zero the matrix column by column? This is as simple as swapping the two for loops. Do this, recompile, and run it. What is going on? Why would simply changing how we loop over the matrix make it take so much longer?

The cause of this requires us to look at the memory hierarchy. There are three main levels of memory:

- Registers: This is where additions and multiplications take place, fast and small
- Cache: Between Registers and Main Memory, intermediate speed and size
- Main Memory: RAM, large and slow



When the processor requires data it first checks to see if it is in the cache. If it finds the data in the cache then this is known as a cache hit. If it doesn't find the data then this is a cache miss and the processor must go out to main memory to get the data. It can take ten times longer to get the data on a cache miss than a cache hit.

When the processor gets memory it does not get a single byte. It gets a larger block of memory called a cache line. This is generally 64 bytes. This is the reason for the large difference in speed between the two programs.

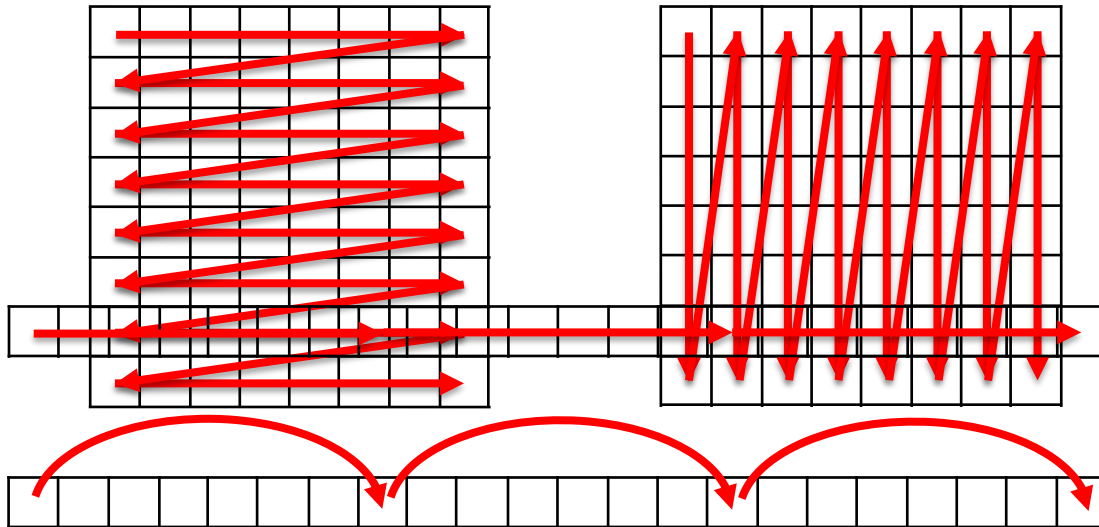
First consider the row by row program. It loads a cache line of the array from main memory, a cache miss, and sets the first element to zero. When it goes to set the second element to zero the cache

line is already in the cache, so we get a cache hit. The next few elements will also be cache hits until we iterate far enough that the elements are in a different cache line.

Now consider the column by column program. The first element results in a cache miss like the row by row program. However, the next element will also be a cache miss as the elements are not next to each other in memory. Every single element results in a cache miss when we loop over the columns.

This is the cause of the performance difference. The row by row program only has to wait on a slow cache miss every eight or so elements while the column by column program cache misses on every element.

As programmers, we cannot change how the processor fetches memory, but by changing how we access our data we can make sure the cache and cache lines help us instead of hurting us.



Now let us see how cache lines might affect parallel programming. Consider the following program:

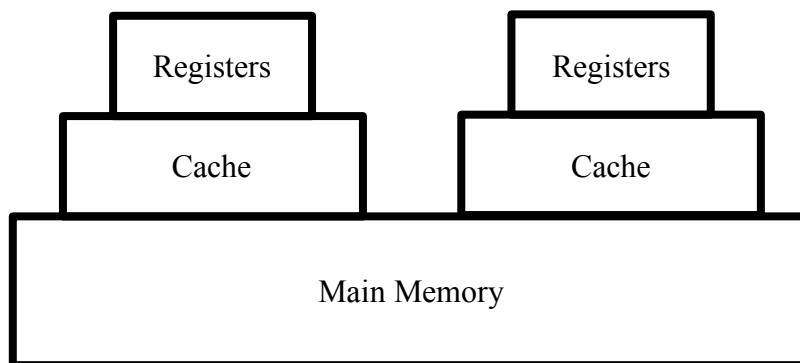
```
1 #include <iostream>
2 #include "sys/time.h"
3 #include "omp.h"
4
5 using namespace std;
6
7 int main()
8 {
9     const int NUM_THREADS = 1;
10    int* array = new int[NUM_THREADS];
11    for(int i = 0; i < NUM_THREADS; i++)
12        array[i]=0;
13    struct timeval tv1, tv2;
14
```

```

15     gettimeofday(&tv1,NULL); //get the current time
16
17     #pragma omp parallel num_threads(NUM_THREADS)
18     {
19         int index = omp_get_thread_num();
20         for(int i = 0; i < 1e9; i++)
21             array[index]++;
22     }
23
24     gettimeofday(&tv2,NULL);
25
26     //this prints out the time in seconds between
27     //the two calls to gettimeofday
28     cout << "Time: " << (double)(tv2.tv_usec - tv1.tv_usec) / 1000000 +
29     (double)(tv2.tv_sec - tv1.tv_sec) << endl;
30
31     delete [] array;
32     return 0;
33 }

```

There are NUM_THREADS threads that each increments their own int in a loop. Try running the program with 1, 2, and more threads. The reason for the decrease in performance is once again due to cache lines. Now this should sound wrong. After all this time array fits on a single cache line. Shouldn't this be very good for performance? The problem is that now we are dealing with multiple threads. The threads will be running on their own cores and each core has its own cache. When the first thread reads array[0] it loads the cache line into the first processor's cache. At the same time the second thread reads array[1] and this loads the same cache line into the second processor's cache. At this point there is not a problem. The same cache line can be in two different caches at the same time. The issue occurs when either thread writes the incremented value back to array. Now the cache line in the other processor's cache is wrong. The other processor must now get the updated cache line from main memory, a cache miss. Now you might think that this is unnecessary in this instance as the first thread does not care about the second thread's array value and vice versa. However, remember that memory transfers work at the cache line level. If thread one kept its cache line despite thread two updating its value, then when thread one writes the cache line back to main memory it would overwrite thread two's value with the old value.



This phenomenon is known as false sharing. The way to prevent it is to make sure each thread is working on different cache lines. To this end consider a modification to our previous code:

```
1 #include <iostream>
2 #include "sys/time.h"
3 #include "omp.h"
4
5 using namespace std;
6
7 int main()
8 {
9     const int NUM_THREADS = 2;
10    const int SPACING = 1;
11    int* array = new int[NUM_THREADS*SPACING];
12    for(int i = 0; i < NUM_THREADS; i++)
13        array[i*SPACING]=0;
14    struct timeval tv1, tv2;
15
16    gettimeofday(&tv1,NULL); //get the current time
17
18    #pragma omp parallel num_threads(NUM_THREADS)
19    {
20        int index = omp_get_thread_num()*SPACING;
21        for(int i = 0; i < 1e9; i++)
22            array[index]++;
23    }
24
25    gettimeofday(&tv2,NULL);
26
27    //this prints out the time in seconds between
28    //the two calls to gettimeofday
29    cout << "Time: " << (double)(tv2.tv_usec - tv1.tv_usec) / 1000000 +
30    (double)(tv2.tv_sec - tv1.tv_sec) << endl;
31
32    delete [] array;
33    return 0;
34 }
```

Now the integers in `array` do not have to be next to each other. Instead they have indexes that are a multiple of `SPACING`. Play around with different `SPACING` sizes; once it gets big enough the slow performance should go away as each value is on its own cache line.