

**E0 – 243**  
**Computer Architecture**  
**Assignment 3**  
**(IISc Bangalore)**

**GitHub Repository**  
**[https://github.com/karm-patel/Assignment\\_3](https://github.com/karm-patel/Assignment_3)**

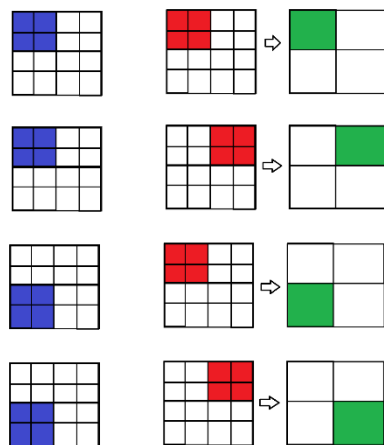
**By**  
**Yash Punjabi, Karm Patel**  
**[yashpunjabi@iisc.ac.in](mailto:yashpunjabi@iisc.ac.in), 20870)**  
**[karmpatel@iisc.ac.in](mailto:karmpatel@iisc.ac.in), 21542)**

# REDUCED MATRIX MULTIPLICATION

## 1. Introduction

We are given implementation of modified matrix multiplication. It takes two ( $N \times N$ ) matrices and gives an output matrix of dimension ( $N/2 \times N/2$ ). Here we have  $N = 2^k$  where  $k \geq 2$ .

It works by taking the first two rows of A and first two columns of B, multiplying them and adding all the elements and placing it into the final output matrix as one entry. This process is continued for all possible pairs of rows and columns. The illustration of above is represented below



The image shows the computation of partial answers in iteration, the same process would be followed. The movement of red block depicts inner loop iteration

Since it is a normal matrix multiplication just with reduced size and also having different blocks multiplication at every iteration which are independent of each other, there can be a potential to optimize the operations and improve the overall performance

Our task is to optimize this naive implementation and create another multi-threaded version as well

## 2. Processor Specifications

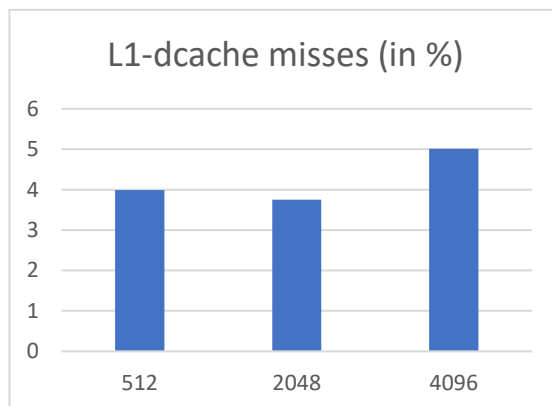
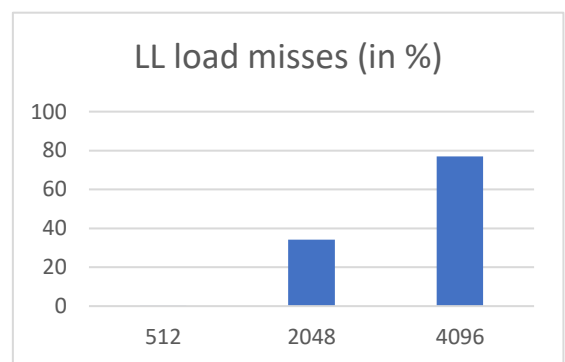
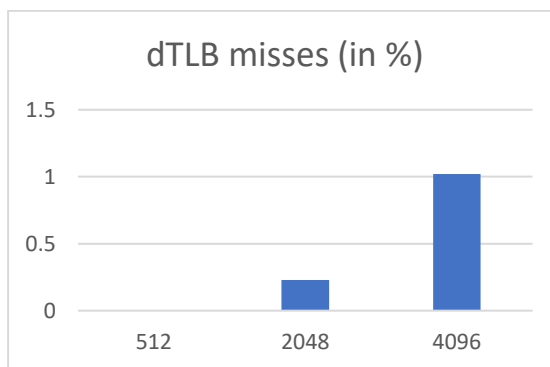
|                  |  |
|------------------|--|
| <b>CPU</b>       | 11th Gen Intel(R) Core (TM) i5-11500 @ 2.70GHz                             |
| <b>Memory</b>    | 16 GB DDR4   |
| <b>Cache</b>     | L1d cache = 48 KB, L1i cache = 32KB,<br>L2 cache = 512KB, L3 cache = 12 MB |
| <b>OS</b>        | Linux Ubuntu 22.04   |
| <b>PAGE SIZE</b> | 4096B  |

### 3. Analysis

We will analyze the above implementation on various matrix sizes in order to identify the ways to optimize and reason about it to improve the overall performance. For analysis perf tool was used. The analysis was done on different input sizes varying from  $N = 16$  to  $N = 16384$ , but for report we will show analysis on three input sizes,  $N = 512, 2048, 4096$ .

#### 3.1 Issues with reduced multiplication and need for blocking

In reduced matrix multiplication for each two rows of matrix A, matrix B two columns are accessed in column wise. But we know matrices are stored in memory in either row major or column order, since here we are trying to access matrix in both fashions, we can see the amount of miss percentage in caches.



By seeing the above graphs, it can be observed reduced matrix multiplication has high LL miss rate and L1-dcache miss rate. It is because of high number of stalls. In LL miss for 2048 we have miss rate of roughly 34% and we know accessing memory is costly, further increasing size of  $N$  will increase this value.

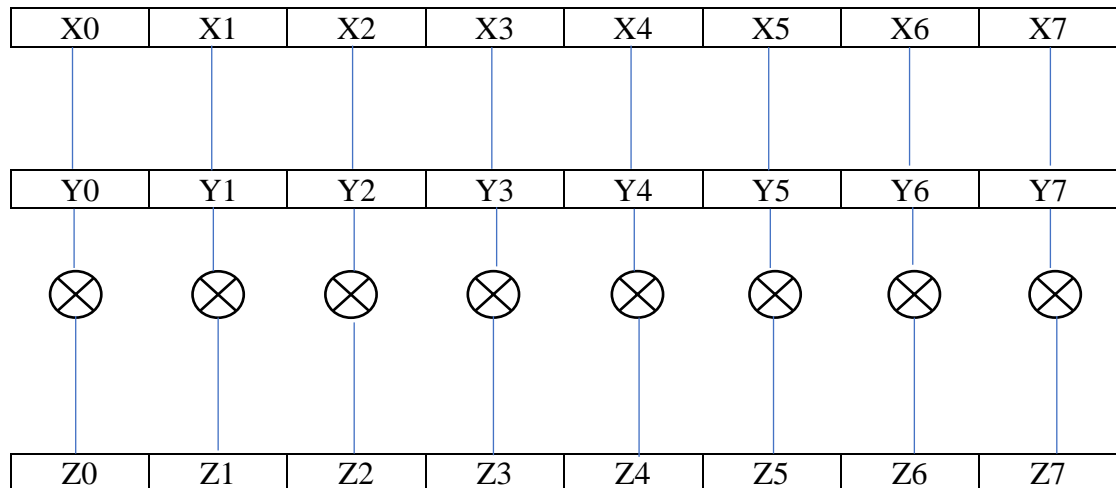
Note that dTLB miss rate is there but has low value. For 512 we have around 0% of miss. This can be probably because  $N$  is not large enough to capture the result.

### 3.2 Optimizing the reduced matrix multiplication

In the reduced method, elements of A were accessed in row major order and whereas elements of matrix B are accessed in column major. Also, the computation for one block to another block of matrix B is independent. Therefore, we can perform multiplicative and additive operation in parallel.

#### 3.2.1 Single Thread optimization

- a. Use of Vector instructions (AVX2 of size 256 bits): For performing operations in parallel, since integer are 32 bits, we can perform 8 32-bit values in one operation



- b. Ijk loop was used in order to overcome access efficiency

Approach:

First two elements of 1<sup>st</sup> row of matrix A were loaded in AVX-256 register r1, r2. Similarly next two elements of 2<sup>nd</sup> row of matrix A were loaded in AVX-256 register r3, r4. Then first two rows of matrix B were loaded on AVX register col1 and col2, following operation were performed

$$M1 = R1 * col1,$$

$$M2 = R3 * col1,$$

$$M3 = R2 * col2,$$

$$M4 = R4 * col2$$

Addition operation was performed along M1 and M2 giving output as AVX register of 8 values. The adjacent values of resultant array were added giving values of output index

The image shows a 15x15 grid. The top-left 2x2 area is colored: blue (row 1, col 1), yellow (row 1, col 2), green (row 2, col 1), and grey (row 2, col 2). The bottom 5 rows (rows 11-15) are shaded light blue. The remaining grid cells are white.

In first iteration AVX register

R1 contains all 8 values of shown color blue

R2 contains all 8 values of shown  
color yellow

R3 contains all 8 values of shown color green

R4 contains all 8 values of shown color grey

The values will change when column iteration completes

In iteration 1:

(Color shown in red) is loaded in col1 AVX register

(Color shown in dark grey) is loaded  
in col2 AVX register

Thus, after inner loop column iteration column will be incremented by 8 and both col1 and dcol2 will slide over while r1, r2, r3, r4 values will remain same in that iteration.

A 10x10 grid with the top-left 8x2 area shaded. The top row of the shaded area is red, and the bottom row is dark blue.

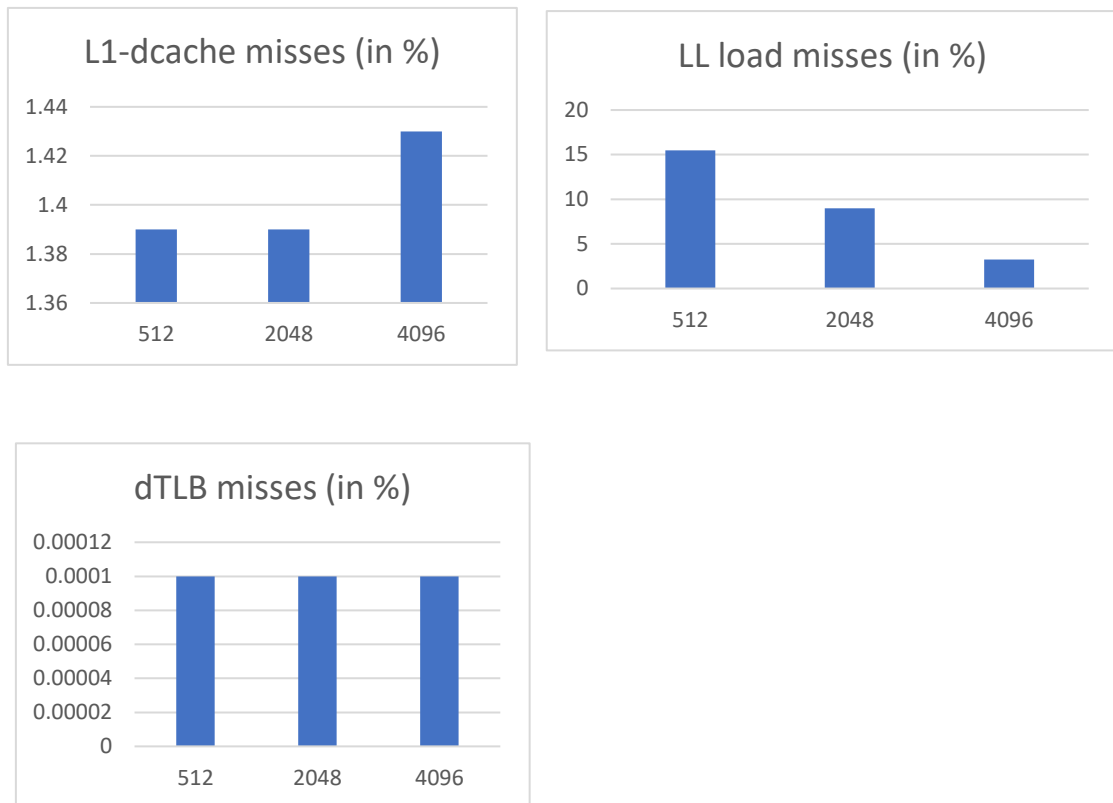
**(MAT B)**

**(Example for N=16)**

Since we have performed horizontal operations the values for output array would be partial values and would be updated in every iteration

|         |         |         |         |         |         |         |         |
|---------|---------|---------|---------|---------|---------|---------|---------|
| Sum [0] | Sum [1] | Sum [2] | Sum [3] | Sum [4] | Sum [5] | Sum [6] | Sum [7] |
|---------|---------|---------|---------|---------|---------|---------|---------|

### Cache Misses



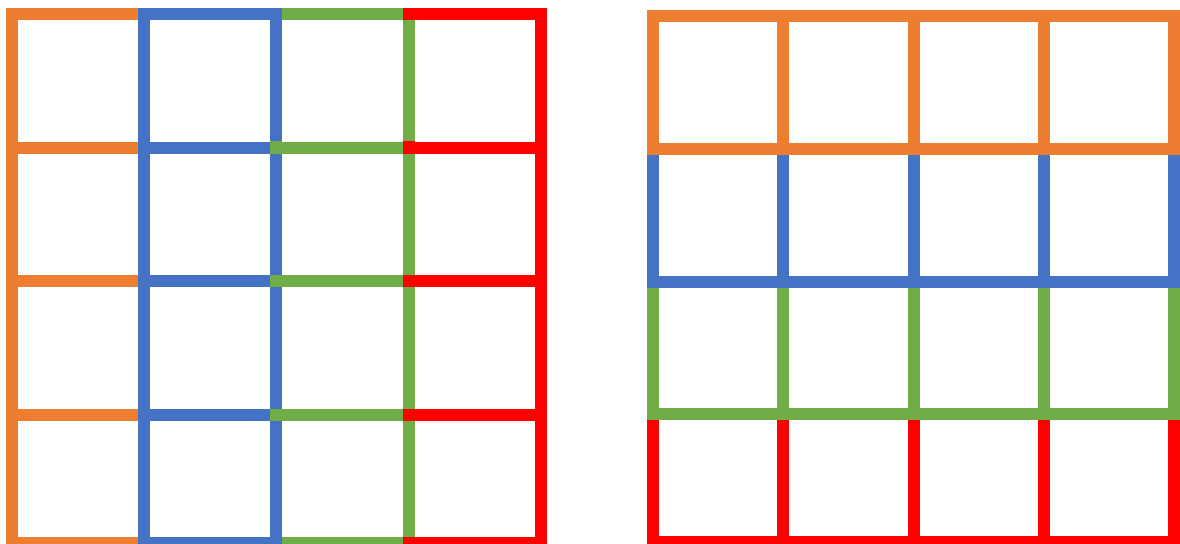
The above results show margin of values between l1-dcache misses and LL load misses. Thereby telling misses have decreased improving the overall performance (The execution time comparison shows the same).

We can see for  $N = 4096$  the LL miss rate from ~77% has reduced to 3%, telling majority of load access would have been found in the L1-dcache itself.

### 3.3 Multi Thread Optimization

Multithreading is the method to run different parts of program in parallel across different cores, ensuring parallelism. It does through using the help of threads which are nothing but lightweight process with the process. A thread can run on different cores leading to maximum utilization of CPU. The advantages of using multithreading

- a. Resource Sharing: There are two ways for process to share resource either Message passing or shared memory. This has to be done explicitly done by programmer. Therefore, the benefit of sharing data allows process to have several threads within same address space
- b. Scalability: The benefits of multithreading excel in multiprocessor architecture. Multi-threading on multiple cores increases the parallelism
- c. Responsiveness: Let's say multithreading is not implemented on a system which receives a request performs a set of task and returns it sequentially. If long process comes system would be busy doing that and make other requests wait. While in case of multithreaded system may allow run a program even if a part of is blocked or doing length operation, thus increasing responsiveness



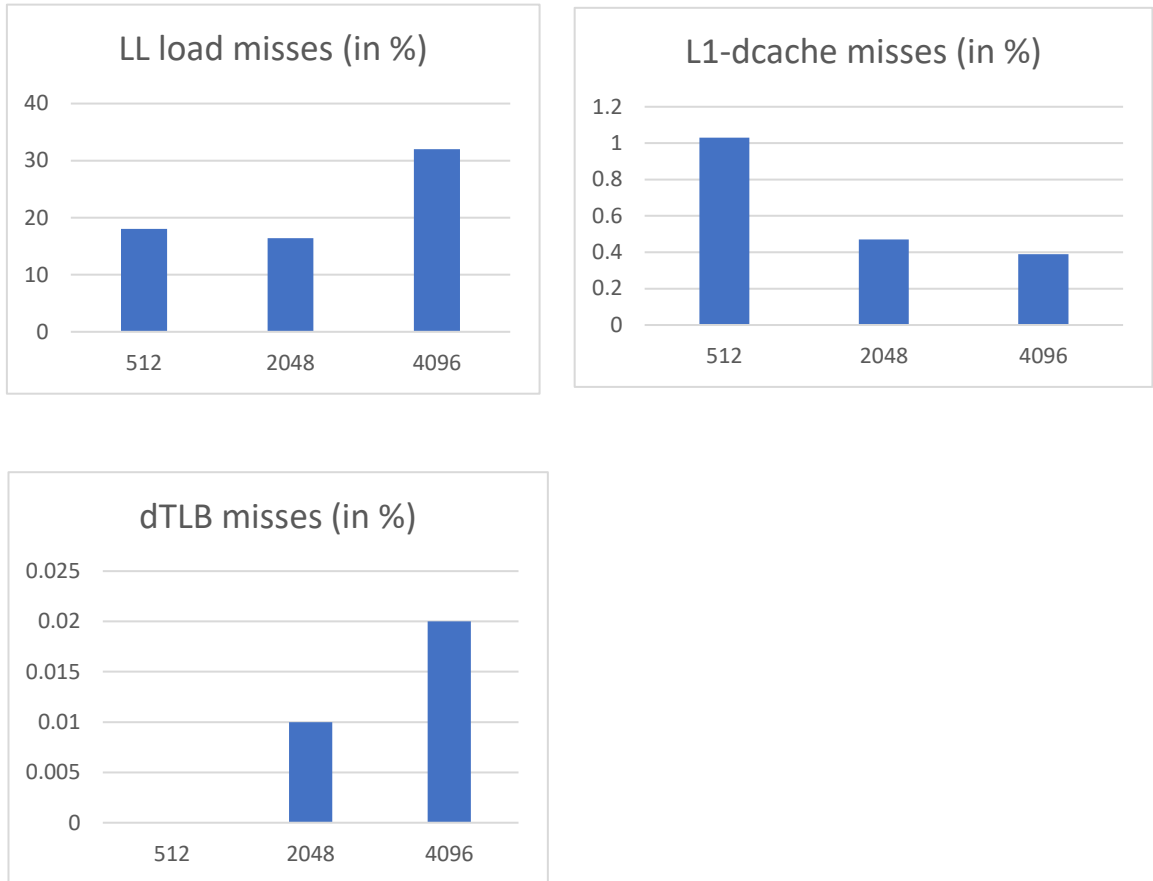
Approach:

Here our program can take advantage of private caches of each core. Therefore, instead of depending on single L1 and L2 cache as in case of above we have multiple now

Total 64 threads are created but they are created in iteration of 4, that is  $64/4 = 16$ , 16 threads in an iteration which computes the value of one color shown above and that is executed 4 times. Since data is independent parallelism can be achieved and there will be no need of mutex locks, since no race condition would be present

Above we have created 4 blocks in row and col that means matrix of size N is divided into  $N/4$  rows and  $N/4$  cols. Since the multiplication of blocks within thread (i.e  $N/4 * N/4$ ) can be further optimized it is performed using tiling along with the vectorized approach which we discussed above in single thread

### 3.3.1 Results



The above graph shows the decrease in miss rate compared to single thread and in reference. In case of 4096 we have ~31% miss rate and in the case of L1-d cache miss rate has decline along with increase in size (for e.g., 4096 has ~ 0.4% miss rate)

dTLB has low percentage it can be due to fact that size N is not large enough to capture the details

The system on which analysis was done had 12 cores. (For size  $\leq 128$  multithread was not used, since it did not provide much of performance boost compared to original version)



### 3.4 Conclusion

| Matrix Size | Reduced Matrix (Naive version) (MS) | Single Thread (MS) | Multithread (MS) |
|-------------|-------------------------------------|--------------------|------------------|
| 16          | 0.025                               | 0.025              | 0.007            |
| 128         | 3.134                               | 0.988              | 3.776            |
| 256         | 24.762                              | 7.586              | 4.359            |
| 512         | 230.973                             | 59.933             | 22.515           |
| 1024        | 1956.72                             | 485.665            | 178.589          |
| 2048        | 44569.6                             | 3797.07            | 1300.02          |
| 4096        | 481741                              | 29931              | 9720.1           |
| 8192        | 5e+06                               | 239245             | 128314           |
| 16384       | 5.76e+7                             | 2.07e+06           | 682936           |

Approx Speedup achieved:

- for single Thread with respect to naive version = 27x (for size 16384)
- for multi-Thread with respect to naive version = 84x (for size 16384)
- for multi-Thread with respect to single Thread = 3x (for size 16384)

An ideal implementation would have achieved at least speed up of 4x when cores are increased but our implementation was not able to achieve this speedup possibly because of the kernel process running in background