

# Assignment 3: Deep-Diving into the Architectural Simulator and Code Optimization

## CPEN 411 – Computer Architecture

Isabelle André – 12521589

11 November 2022

---

### 1.0 Project Description

This assignment consisted in understanding the implementation of architectural simulators and how to generate traces for the simulator from applications. A naive matrix multiplication function was provided to be optimized in speedup and memory. In Part A, we explore the simulator's pipelined architecture and implement a non-pipelined core by modifying the current out-of-order core implementation of ChampSim. In Part B and this report, we optimize the basic matrix multiplication implementation by achieving 30% speed-up using blocked matrix multiplication, and reducing the L1D cache miss-rate by 8x by using loop reordering and matrix transposition.

### 2.0 Matrix Multiplication Optimization

Matrix multiplication is an important operation used in multiple programming applications and convolutions in neural network instances. In applications requiring highly optimized code for machine learning applications, it is beneficial to optimize this operation to have cache-friendly access and avoid data dependencies by restructuring the code, while retaining its initial function.

#### 2.1 Baseline Matrix Multiplication

##### 2.1.1 Description

The naive matrix multiplication implementation uses a triple nested for loop structure to multiply two matrices:

$$C_{ij} = \sum_{k=1}^b A_{ik} * B_{kj}$$

As this method takes an iterative and exhaustive approach to multiplication, this method can be optimized to avoid RAW data dependencies and fill the cache in more effective ways.

##### 2.1.2 Implementation

With a matrix size of N, the total time complexity for this method is  $O(n^3)$ .

```
for(unsigned long long int i=0; i<matrix_size; i++){
```

```

        for(unsigned long long int j=0; j<matrix_size; j++){
            for(unsigned long long int k=0; k<matrix_size; k++){
                mat_c[j*matrix_size + i] += mat_a[k*matrix_size +
i]
                * mat_b[j*matrix_size + k];
            }
        }
    }
}

```

Table 2.1 shows the Naive matrix multiplication performance metrics as a benchmark for comparison, including cumulative IPC, Cache miss rates, and Branch prediction accuracy.

	Baseline
<b>Cumulative IPC</b>	0.756089
<b>L1D Miss Rate</b>	266,451.00
<b>L2C Miss Rate</b>	137,886.00
<b>LLC Miss Rate</b>	5,587.00
<b>Branch Prediction Accuracy</b>	99.2247

Table 2.1: Naive Matrix Multiplication Baseline

## 2.2 Matrix Multiplication Optimizations

As most algorithms process and execute operations over large quantities of data, optimizing the cache access and usage is the first and most efficient way to improve performance. Memory acts as a bottleneck to CPU performance, and by improving memory usage, it is also possible to improve other performance metrics such as IPC.

### 2.2.1 Loop Reordering

Loop reordering explores the effects of locality of reference in cache memory. As the matrix is stored in row-major order, everytime the inner-loop is incremented, an entire row of the matrix is skipped, therefore jumping further into memory blocks than cached values. Applying reordering improves performance gains by increasing the spatial locality of data. This presents an optimization allowing to eliminate cache misses on each iteration of the inner loop i.

### 2.2.2 Loop Reordering Implementation

This method increases the time complexity by  $O(n^2)$  as compared to the previous.

```

for(unsigned long long int j=0; j<matrix_size; j++){
    for(unsigned long long int k=0; k<matrix_size; k++){
        for(unsigned long long int i=0; i<matrix_size; i++){

```

```

        mat_c[j*matrix_size + i] += mat_a[k*matrix_size +
i]
        * mat_b[j*matrix_size + k];
    }
}
}

```

### 2.2.3 Matrix Transpose Optimization

One of the fundamental problems of naive matrix multiplication is that when iterating through matrix B,  $n$  elements are jumped over as a column is iterated. The Matrix Transpose method consists in storing matrix B and transposing it prior to the matrix multiplication. While this optimization increases the time complexity by  $O(n^2)$ , sequential reads are guaranteed in the 3rd nested loop, optimizing serial caching access and spatial locality. Spatial locality may be even further optimized when used in conjunction with loop reordering or blocking. In our code, matrix transpose is used in combination with loop reordering previously used in Section 2.2 to further optimize it.

This method is implemented in **matmul\_opt2** of Part B of this assignment, to optimize the memory usage and decrease L1D cache misses by over 8x.

### 2.2.4 Matrix Transpose Optimization Implementation

This method increases the time complexity by  $O(n^2)$  as compared to the previous.

```

unsigned long long int new_mat_b [matrix_size * matrix_size];
// TRANSPOSE
for (int i=0; i < matrix_size; i++) {
    for (int j=0; j < matrix_size; j++) {
        new_mat_b[j*matrix_size + i] = mat_b[i*matrix_size + j];
    }
}

// REORDER
for (int j = 0; j < matrix_size; j++) {
    for (int k = 0; k < matrix_size; k++) {
        for (int i = 0; i < matrix_size; i++) {
            mat_c[j*matrix_size + i] += mat_a[k*matrix_size + i]
            * mat_b[j*matrix_size + k];
        }
    }
}
}

```

### 2.2.5 Memory Blocking Optimization

Memory blocking can be used to optimize the caching mechanism, employing temporal locality properties of caches. In block matrix multiplication, the matrix is

divided into multiple  $S \times S$  sub-matrices before multiplying corresponding blocks. In our implementation,  $S$  is equivalent to cache block size of 64. Blocking allows the code to store a block into the cache and complete all writes and reads at once before evicting the block and repeating the process for the next, therefore improving the temporal locality of the nested for loops. This method is used in addition to loop reordering to further optimize both IPC and caching efficiency.

We use this method in **matmul\_opt 1** of Part B of this assignment to optimize the IPC by 30%, and while a noticeable improvement is seen in L2C and LLC cache misses, there is little improvement shown for the L1D cache.

### 2.2.6 Memory Blocking Implementation

```
for (unsigned long long int i = 0; i < matrix_size; i +=
BLOCK_SIZE) {
    for (unsigned long long int j = 0; j < matrix_size; j +=
        BLOCK_SIZE) {
        mat_c[i*matrix_size+j] = 0;
        for (unsigned long long int k = 0; k < matrix_size; k
            += BLOCK_SIZE) {
            for (unsigned long long int x = i; x < ((i +
                BLOCK_SIZE) > matrix_size ? matrix_size : (i
                    +
                        BLOCK_SIZE)); x++) {
                for (unsigned long long int y = j; y < ((j +
                    BLOCK_SIZE) > matrix_size ? matrix_size : (j
                        + BLOCK_SIZE)); y++) {
                    for (unsigned long long int z = k; z < ((k
                        + BLOCK_SIZE) > matrix_size ?
                            matrix_size
                                : (k + BLOCK_SIZE)); z++) {
                        mat_c[x*matrix_size + y] +=
                            mat_a[z*matrix_size + y]*mat_b[x*
                                matrix_size + z ];
                    }
                }
            }
        }
    }
}
```

### 2.3 Speedup and Memory Comparison

A breakdown of the performance metrics for the naive matrix multiplication and both of the above optimized matrix multiplications are shown in Table 2.2.

	Baseline	Reorder	Transposition	Blocking
--	----------	---------	---------------	----------

<b>Cumulative IPC</b>	0.756089	1.588630	1.968900	1.07228
<b>L1D Miss Rate</b>	266,451	28,898	24,686	221636
<b>L2C Miss Rate</b>	137,886	15,604	12,503	187
<b>LLC Miss Rate</b>	5,587	5,240	4,950	187
<b>Branch Prediction Accuracy</b>	99.2247	99.2247	99.2231	96.971

Table 2.2: Performance Metrics Breakdown for Naive and Optimized Matrix Multiplication

### 2.2.2 Speedup

The cumulative IPC obtained by Reordering the nested for loops from ijk to a jki configuration showed an improvement in IPC of 1.902020 compared to the previous naive method at 0.756089.

$$Reorder Speedup = \frac{Reorder IPC}{Naive IPC} = \frac{1.588630}{0.756089} = 2.10111$$

Therefore, the Reorder method showed an improvement of 2.10111 as compared to the Naive method.

The cumulative IPC obtained by Transposing matrix B prior to the matrix multiplication showed a large improvement in IPC at 1.9689.

$$Transpose Speedup = \frac{Transpose IPC}{Naive IPC} = \frac{1.9689}{0.756089} = 2.60405$$

Therefore, the Reorder-Transposing method showed an improvement of 2.60405 as compared to the Naive method.

The cumulative IPC obtained using cache blocking to matrix multiply only shows a slight increase in IPC at 1.07228.

$$Blocking Speedup = \frac{Blocking IPC}{Naive IPC} = \frac{1.07228}{0.756089} = 1.41819$$

Therefore, the Blocking method showed an improvement of 1.41819 as compared to the Naive method.

### 2.3.3 Memory Optimization Improvement

$$Reorder L1D Miss Improvement = \frac{Baseline L1D Miss}{Reorder L1D Miss} = \frac{266451}{28898} = 9.22x$$

As opposed to the previous naive multiplication method, the Reorder optimization shows a reduction in L1D cache misses of 9.22x.

$$Transpose L1D Miss Improvement = \frac{Baseline L1D Miss}{Transpose L1D Miss} = \frac{266451}{24686} = 10.8x$$

The Transpose optimization shows the greatest reduction in L1D cache misses of 10.8x.

$$\text{Blocking L1D Miss Improvement} = \frac{\text{Baseline L1D Miss}}{\text{Transpose L1D Miss}} = \frac{266451}{221636} = 1.20x$$

Though the Blocking method shows a great reduction in L2C and LLC cache misses as shown in Table 2.2, it only shows a small cache miss reduction in L1D of 1.20x.

## 2.4 Plotted Comparison

Comparing each method's cumulative IPC, it is observed that while the Reorder method has a large increase in performance, Transposition even further increases its computational density, as further improvements to spatial locality are applied. Figure 2.1 shows the IPC of the optimized matrix multiplications as compared to the baseline naive method.

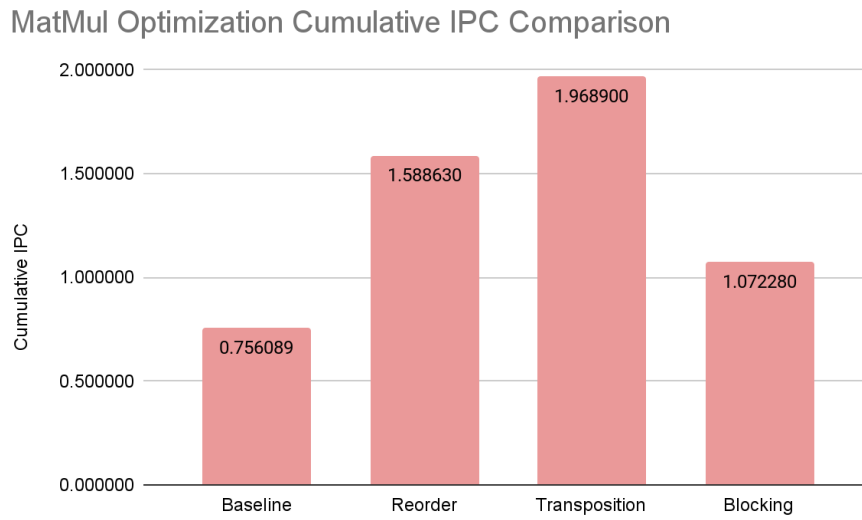


Figure 2.1: IPC Metrics Comparison for Naive and Optimized Matrix Multiplication

Figure 2.2 compares the cache miss rates of the optimized matrix multiplications as compared to the baseline naive method. While the LLC cache miss rates are all approximately the same for Naive, Reordering, and Transpose methods, a large improvement is seen when reordering and transposing for L1D and L2C. Memory blocking does not show much improvement in cache misses in the L1D, but shows a much better performance in the L2C and LLC.

Figure 2.3 shows the branch prediction accuracy of each method, and shows little improvement in any of the optimizations.

### MatMul Optimization Cache Miss Rate Comparison

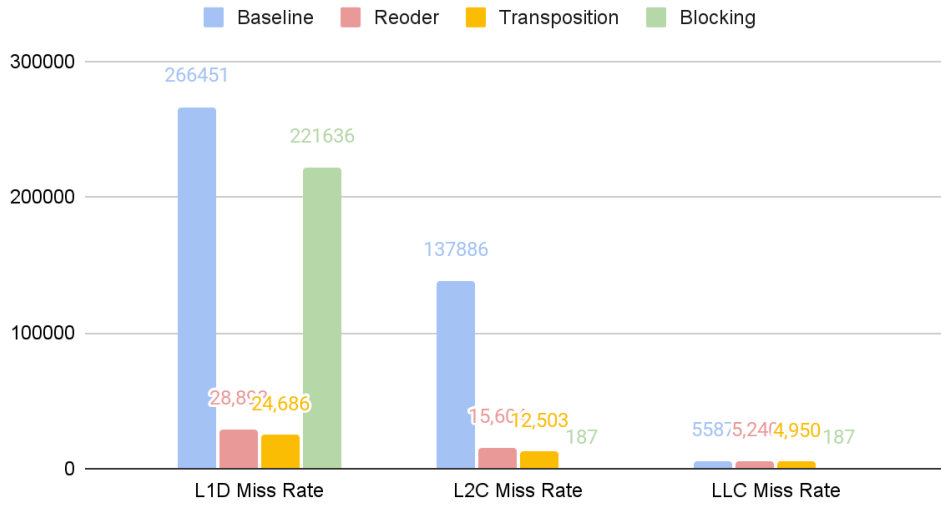


Figure 2.2: Cache Miss Metrics Comparison for Naive and Optimized Matrix Multiplication

### MatMul Optimization Branch Prediction Accuracy Comparison

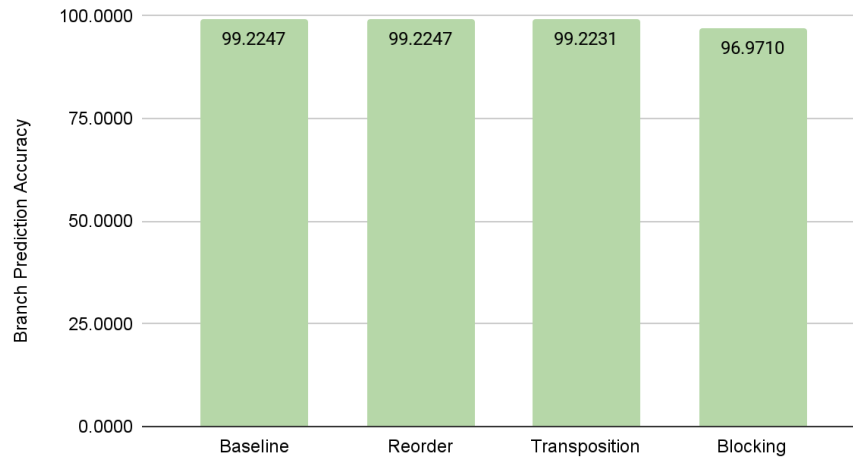


Figure 2.3: Branch Prediction Metrics Comparison for Naive and Optimized Matrix Multiplication

## 4.0 Conclusion

As matrix multiplication is a fundamental operation used in deep learning, its optimization is a topic of interest in many scientific applications. It can be improved in a variety of ways using methods to improve its spatial and temporal locality. **Memory blocking** was used to improve the spatial and temporal locality of the cache, showing a **speedup of 1.4x** that of the naive method. While there is little improvement in L1D cache misses, it shows a much better performance in the L2C and LLC.

Loop reordering was then to improve the spatial locality of the cache accesses, by reordering the nested for loops such that it is less likely to incur cache misses in the inner loop. **Matrix transposition** was then used to store and transpose matrix B prior to multiplication to guarantee sequential reads in the 3rd nested loop. This optimizes serial caching access and spatial locality, **decreasing L1D cache misses by 10.8x**.

There are many more memory optimizations to explore individually or in conjunction with some of the methods presented in this assignment, such as row-wise multiplication, vectorizing, etc.

## 5.0 References

Sarband, N.M., Gustafsson, O. & Garrido, M. Using Transposition to Efficiently Solve Constant Matrix-Vector Multiplication and Sum of Product Problems. J Sign Process Syst 92, 1075–1089 (2020).

<https://doi.org/10.1007/s11265-020-01560-z>

Hennessy, John L & Patterson, David A. (2012). Computer Architecture, A Quantitative Approach. Morgan Kaufmann. 5th Edition.

[http://acs.pub.ro/~cpop/SMPA/Computer%20Architecture%20A%20Quantitative%20Approach%20\(5th%20edition\).pdf](http://acs.pub.ro/~cpop/SMPA/Computer%20Architecture%20A%20Quantitative%20Approach%20(5th%20edition).pdf)