

Benchmark Suite: Matrix-Matrix Multiplication (mmult)

Mohammed Mansour

April 20, 2025

1 Introduction

Matrix-matrix multiplication (mmult) is a fundamental operation in scientific computing, machine learning, and computer architecture. This report documents the implementation and evaluation of the mmult benchmark in the YABMS benchmark suite. The objective is to understand the performance characteristics of different datasets and analyze the impact of various optimization techniques.

2 Collaboration

I discussed the details of the project with Moneer Al-Bokhaiti. We covered the project requirements, the necessary tools, and the complexity of the code. He suggested creating a separate application to generate the dataset. We also agreed that storing the dataset in a binary file would be more efficient than saving it as ASCII text.

Additionally, I attended a meeting conducted by class students, where we discussed various issues related to the project. During this meeting, we broke the implementation into smaller steps to make the development process more manageable.

3 Naïve Implementation

The naïve implementation is a simple three-loop to compute the matrix product:

```
1 for (int i = 0; i < M; i++) {  
2     for (int j = 0; j < P; j++) {  
3         R[i][j] = 0;  
4         for (int k = 0; k < N; k++) {  
5             R[i][j] += A[i][k] * B[k][j];  
6         }  
7     }  
8 }
```

Listing 1: Naïve mmult implementation

This method is straightforward but inefficient due to cache misses and redundant memory accesses. Each element of the matrices is accessed multiple times, leading to poor utilization of the CPU cache. Additionally, the lack of optimization techniques such as loop unrolling results in nonoptimal performance, especially for larger datasets.

4 Optimized Implementation

The optimized implementation applies a blocking strategy to improve spatial and temporal locality. By computing sub-matrices (blocks) of the result matrix, the approach minimizes cache misses by reusing loaded data efficiently.

Blocking refers to dividing the matrices into smaller square blocks or sub-matrices of size $b \times b$. Instead of computing one element of the result matrix at a time, the algorithm computes partial results over these sub-matrices and accumulates them. This drastically reduces the number of memory accesses to data not residing in cache, which is critical for performance.

The implementation chooses a stationary output sub-matrix, looping over the corresponding input sub-matrices from matrices A and B. This ensures that data is reused as much as possible before being evicted from the cache. The block size b is passed as a command-line argument, making the implementation tunable to fit different cache configurations and matrix sizes.

One key insight is that performance benefits become more significant as the matrix size increases. This is because larger matrices experience more cache misses in the naïve version, and blocking helps mitigate this by exploiting spatial locality in the B matrix and temporal locality in the A matrix.

The optimized implementation is as follows:

```
1 void mmult_opt(float **A, float **B, float **R, int M, int N, int P,  
2   int b) {  
3     for (int ii = 0; ii < M; ii += b) {  
4       for (int jj = 0; jj < P; jj += b) {  
5         for (int kk = 0; kk < N; kk += b) {  
6           for (int i = ii; i < fmin(ii + b, M); ++i) {  
7             for (int j = jj; j < fmin(jj + b, P); ++j) {  
8               float val = R[i][j];  
9               for (int k = kk; k < fmin(kk + b, N); ++k) {  
10                  val += A[i][k] * B[k][j];  
11              }  
12              R[i][j] = val;  
13          }  
14      }  
15  }  
16 }  
17 }
```

Listing 2: Blocked mmult Implementation

Initial bugs involved uninitialized values and improper bounds when indexing. Fixes were made by ensuring the matrices were zero-initialized and carefully managing loop ranges. Additional care was taken to validate that block size does not exceed the dimensions of the matrices to avoid segmentation faults.

5 Dataset Generation

The dataset required for the benchmark is generated using a Python script. The script generates random matrices, performs matrix multiplication, and saves the datasets in binary format. The dataset consists of predefined matrix sizes categorized into five sets: testing, small, medium, large, and native.

The following Python script is used for dataset generation:

```
1 class Dataset:
2     def __init__(self, rowsA, colsA, rowsB, colsB, name):
3         self.rowsA = rowsA
4         self.colsA = colsA
5         self.rowsB = rowsB
6         self.colsB = colsB
7         self.name = name
8
9 dataset = [
10     Dataset(16, 12, 12, 8, "testing"),
11     Dataset(121, 180, 180, 115, "small"),
12     Dataset(550, 620, 620, 480, "medium"),
13     Dataset(962, 1012, 1012, 1221, "large"),
14     Dataset(2500, 3000, 3000, 2100, "native")
15 ]
16
17 def generate_matrix(rows, cols):
18     return np.random.rand(rows, cols).astype(np.float32)
19
20 def save_matrix_binary(filename, matrix):
21     if not os.path.exists("mmult_ds"):
22         os.makedirs("mmult_ds")
23     with open("mmult_ds/" + filename, 'wb') as f:
24         f.write(matrix.tobytes())
25
26 def main():
27     parser = argparse.ArgumentParser(description='Generate and save
28 matrices in binary format.')
29     parser.add_argument('name', type=str, help='Dataset name')
30     args = parser.parse_args()
31
32     dataset_info = next((d for d in dataset if d.name == args.name),
33 None)
34     if not dataset_info:
35         print(f"No dataset found with the name '{args.name}'.")
36         sys.exit(1)
37
38     rowsA, colsA, rowsB, colsB = dataset_info.rowsA, dataset_info.colsA
39     , dataset_info.rowsB, dataset_info.colsB
40     matrixA = generate_matrix(rowsA, colsA)
41     matrixB = generate_matrix(rowsB, colsB)
42     result_matrix = np.dot(matrixA, matrixB)
43
44     save_matrix_binary("matrixA.bin", matrixA)
45     save_matrix_binary("matrixB.bin", matrixB)
46     save_matrix_binary("matrixC.bin", result_matrix)
47
48     print("Matrices and metadata saved successfully in binary format.")
```

Listing 3: Dataset Generation Script

This script ensures efficient dataset generation by leveraging NumPy for matrix operations and storing the results in a structured binary format. The metadata file provides dimensions for easy retrieval and compatibility with benchmarking implementations.

6 Evaluation

The performance of the naïve implementation was evaluated using various dataset sizes. Execution time was measured for each case, and results were compared with optimized implementations.

6.1 Experimental Setup

The experiments were conducted on **Personal Laptop** with the following setup:

- Runner Type: Ubuntu-24.04.1 LTS
- Operating System: Linux
- Processor (CPU): AMD PRO A12-9800B R7, 12 COMPUTE CORES 4C+8G
- Memory (RAM): 8 GB
- Storage (SSD): 512 GB
- Architecture: x64
- Compiler: GCC 13.3.0 with optimization flags -O3

6.2 Results

Table 1 summarizes the execution times for different dataset sizes.

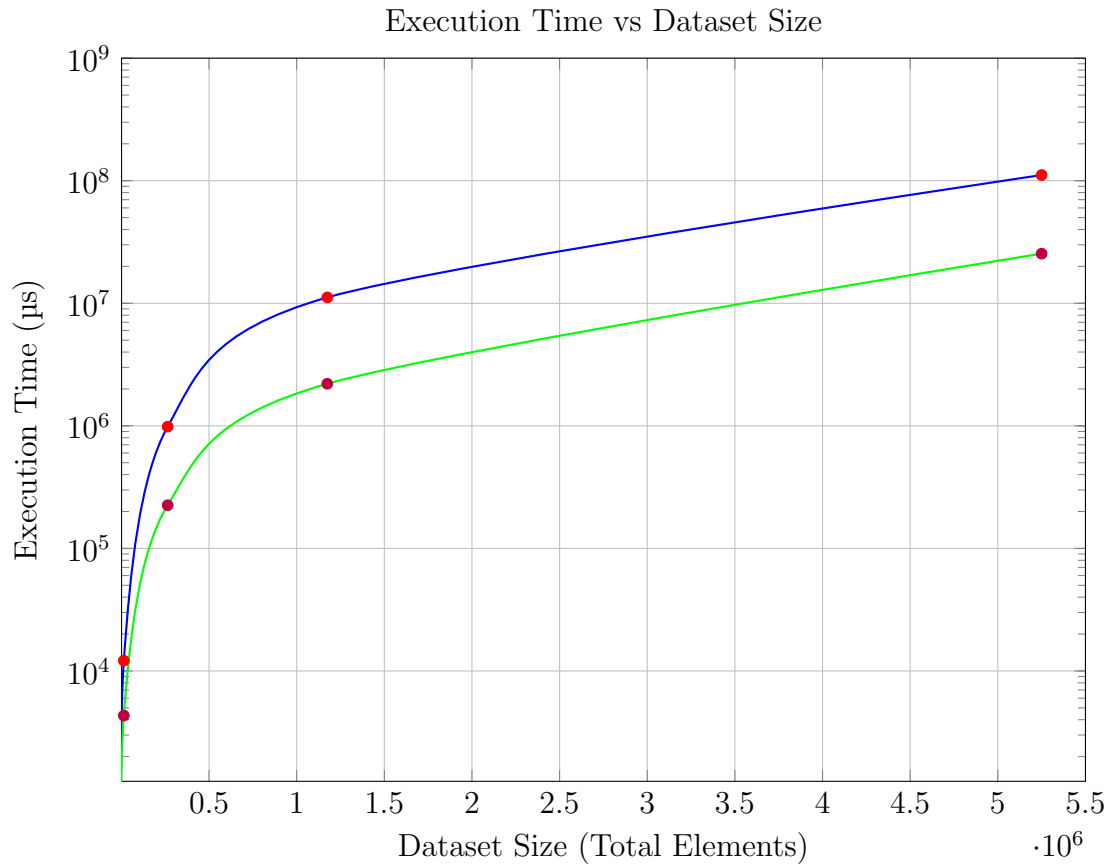
Dataset Size	Naïve Execution Time(ms)	Opt. Execution Time(ms)
Testing (16x8)	0.005	0.003
Small (121x115)	12.137	4.326
Medium (550x480)	984.68	225
Large (962x1221)	11167	2204
Native (2500x2100)	111219	25396

Table 1: Execution times for different dataset sizes.

The benchmark execution logs and results can be accessed at the following link:
https://github.com/mohammed0x00/YABMS/tree/main/docs/runtime_data_opt.

6.3 Graphical Representation

The following graph illustrates the execution time versus dataset size for the naïve implementation:



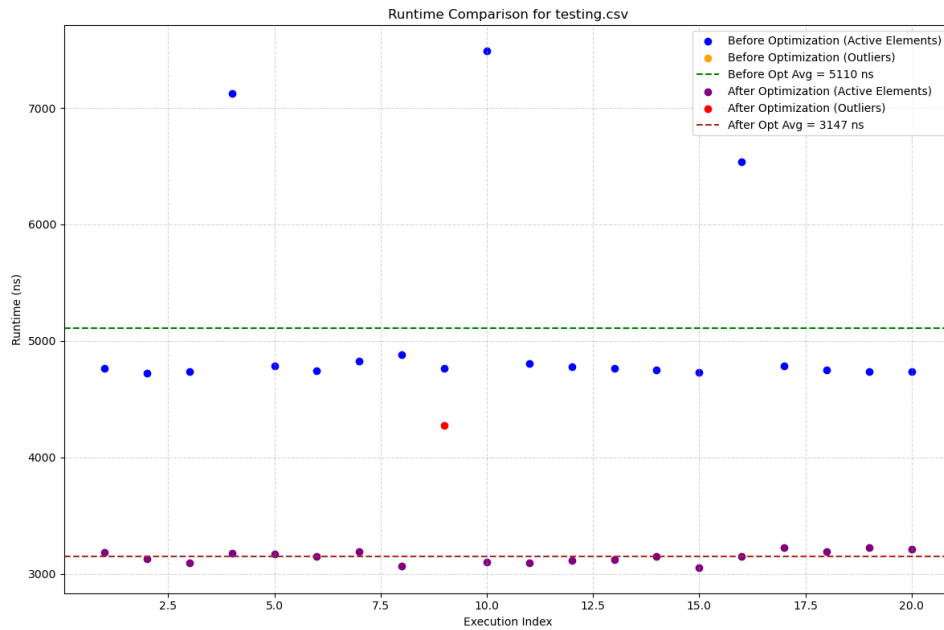
6.4 Detailed Results

The execution times for different dataset sizes, measured in milliseconds (ms), along with their standard deviations. As dataset size increases, execution time rises significantly. The "Testing" dataset (16×8), has an execution time of only 0.005 ms with a negligible standard deviation of 0.0008. The "Small" dataset (121×115) takes 12.137 ms, with a slightly higher deviation of 1.189. The "Medium" dataset (550×480) requires around 1 second, showing a noticeable jump in both execution time and deviation (20.159). The "Large" dataset (962×1221) exhibits a much higher execution time of 11 seconds with a deviation of 46.958. Finally, the "Native" dataset (2500×2100) has the highest execution time of 111 seconds and the highest standard deviation of 661.744, indicating increased variability in execution time as dataset size grows.

The detailed results for each dataset size are presented below:

6.5 Testing Dataset Details

The following are the detailed results for the Testing dataset:

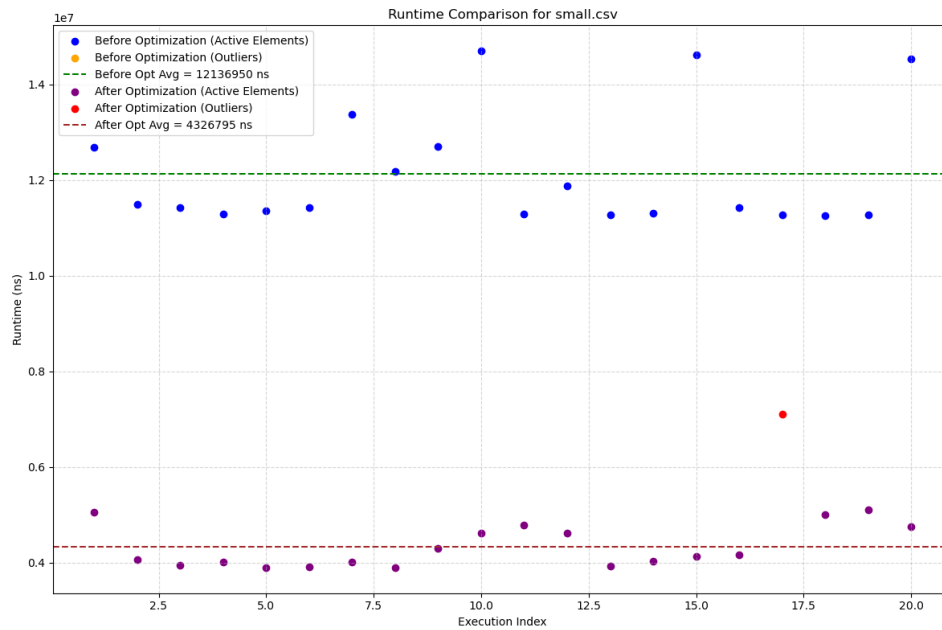


```
1 Setting up schedulers and affinity:
2   * Setting the niceness level:
3     -> trying niceness level = -20
4     + Process has niceness level = -20
5   * Setting up FIFO scheduling scheme and high priority ... Succeeded
6   * Setting up scheduling affinity ... Succeeded
7
8 Running "scalar_opt" implementation:
9   * Invoking the implementation 20 times .... Finished
10  * Verifying results .... Success
11  * Running statistics:
12    + Starting statistics run number #1:
13      - Standard deviation = 250
14      - Average = 3203
15      - Number of active elements = 20
16      - Number of masked-off = 1
17    + Starting statistics run number #2:
18      - Standard deviation = 50
19      - Average = 3147
20      - Number of active elements = 19
21      - Number of masked-off = 0
22  * Runtimes (MATCHING): 3147 ns
23  * Dumping runtime informations:
24    - Filename: scalar_opt_runtimes.csv
25    - Opening file .... Succeeded
26    - Writing runtimes ... Finished
27    - Closing file handle .... Finished
```

Listing 4: Testing Dataset Output

6.6 Small Dataset Details

The following are the detailed results for the Small dataset:

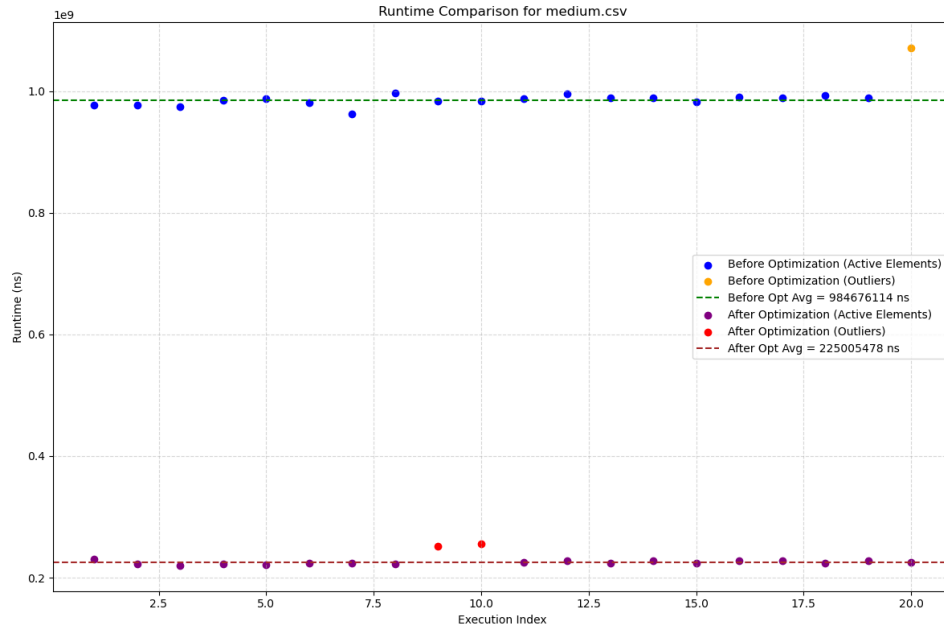


```
1 Setting up schedulers and affinity:
2   * Setting the niceness level:
3     -> trying niceness level = -20
4     + Process has niceness level = -20
5   * Setting up FIFO scheduling scheme and high priority ... Succeeded
6   * Setting up scheduling affinity ... Succeeded
7
8 Running "scalar_opt" implementation:
9   * Invoking the implementation 20 times .... Finished
10  * Verifying results .... Success
11  * Running statistics:
12    + Starting statistics run number #1:
13      - Standard deviation = 733060
14      - Average = 4466056
15      - Number of active elements = 20
16      - Number of masked-off = 1
17    + Starting statistics run number #2:
18      - Standard deviation = 421649
19      - Average = 4326795
20      - Number of active elements = 19
21      - Number of masked-off = 0
22  * Runtimes (MATCHING): 4326795 ns
23  * Dumping runtime informations:
24    - Filename: scalar_opt_runtimes.csv
25    - Opening file .... Succeeded
26    - Writing runtimes ... Finished
27    - Closing file handle .... Finished
```

Listing 5: Small Dataset Output

6.7 Medium Dataset Details

The following are the detailed results for the Medium dataset:



```

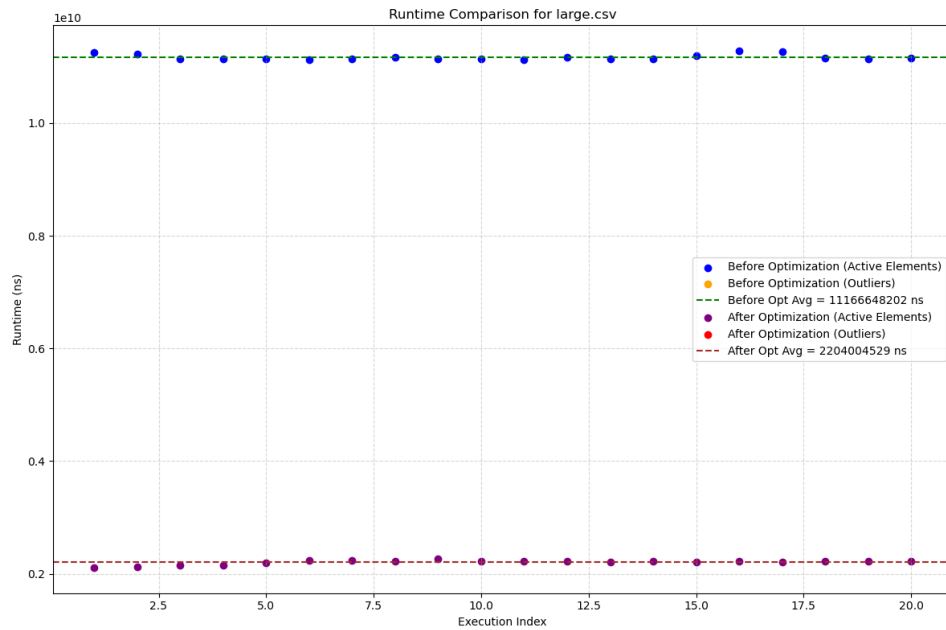
1 Setting up schedulers and affinity:
2   * Setting the niceness level:
3     -> trying niceness level = -20
4     + Process has niceness level = -20
5   * Setting up FIFO scheduling scheme and high priority ... Succeeded
6   * Setting up scheduling affinity ... Succeeded
7
8 Running "scalar_opt" implementation:
9   * Invoking the implementation 20 times .... Finished
10  * Verifying results .... Success
11  * Running statistics:
12    + Starting statistics run number #1:
13      - Standard deviation = 8906230
14      - Average = 227847343
15      - Number of active elements = 20
16      - Number of masked-off = 1
17    + Starting statistics run number #2:
18      - Standard deviation = 6446166
19      - Average = 226399194
20      - Number of active elements = 19
21      - Number of masked-off = 1
22    + Starting statistics run number #3:
23      - Standard deviation = 2637246
24      - Average = 225005478
25      - Number of active elements = 18
26      - Number of masked-off = 0
27  * Runtimes (MATCHING): 225005478 ns
28  * Dumping runtime informations:
29    - Filename: scalar_opt_runtimes.csv
30    - Opening file .... Succeeded
31    - Writing runtimes ... Finished
32    - Closing file handle .... Finished

```

Listing 6: Medium Dataset Output

6.8 Large Dataset Details

The following are the detailed results for the Large dataset:

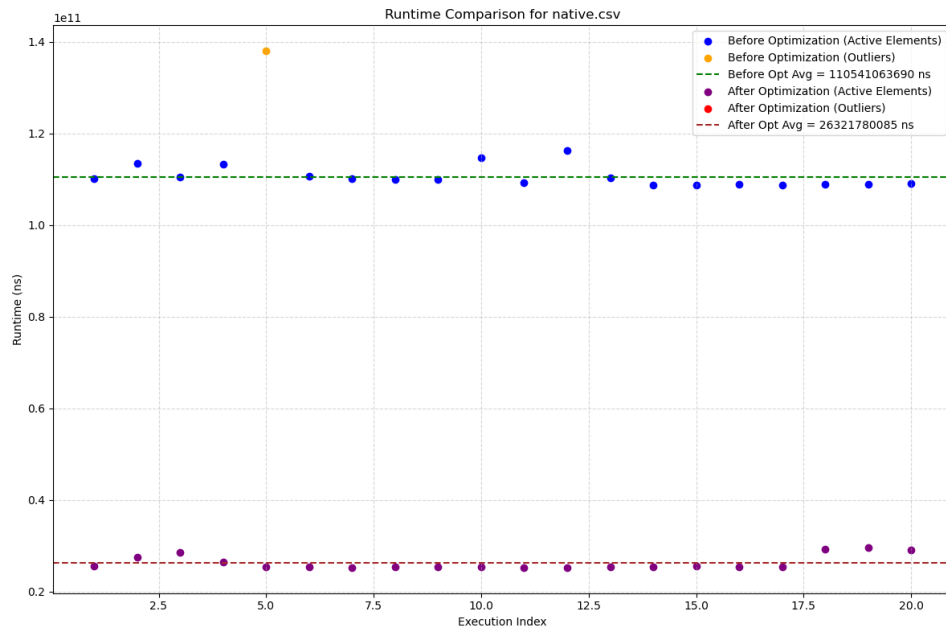


```
1 Setting up schedulers and affinity:
2   * Setting the niceness level:
3     -> trying niceness level = -20
4     + Process has niceness level = -20
5   * Setting up FIFO scheduling scheme and high priority ... Succeeded
6   * Setting up scheduling affinity ... Succeeded
7
8 Running "scalar_opt" implementation:
9   * Invoking the implementation 20 times .... Finished
10  * Verifying results .... Success
11  * Running statistics:
12    + Starting statistics run number #1:
13      - Standard deviation = 38622316
14      - Average = 2204004529
15      - Number of active elements = 20
16      - Number of masked-off = 0
17  * Runtimes (MATCHING): 2204004529 ns
18  * Dumping runtime informations:
19    - Filename: scalar_opt_runtimes.csv
20    - Opening file .... Succeeded
21    - Writing runtimes ... Finished
22    - Closing file handle .... Finished
```

Listing 7: Large Dataset Output

6.9 Native Dataset Details

The following are the detailed results for the Native dataset:



```

1 Running "scalar_opt" implementation:
2 * Invoking the implementation 20 times .... Finished
3 * Verifying results .... Success
4 * Running statistics:
5   + Starting statistics run number #1:
6     - Standard deviation = 734113245
7     - Average = 26321780085
8     - Number of active elements = 20
9     - Number of masked-off = 4
10  + Starting statistics run number #2:
11    - Standard deviation = 574367356
12    - Average = 25598762400
13    - Number of active elements = 16
14    - Number of masked-off = 1
15  + Starting statistics run number #3:
16    - Standard deviation = 286857123
17    - Average = 25468953877
18    - Number of active elements = 15
19    - Number of masked-off = 1
20  + Starting statistics run number #4:
21    - Standard deviation = 94935667
22    - Average = 25396312350
23    - Number of active elements = 14
24    - Number of masked-off = 0
25 * Runtimes (MATCHING): 25396312350 ns
26 * Dumping runtime informations:
27   - Filename: scalar_opt_runtimes.csv
28   - Opening file .... Succeeded
29   - Writing runtimes ... Finished
30   - Closing file handle .... Finished

```

Listing 8: Native Dataset Output

7 Conclusion

The benchmarking results demonstrate significant differences in execution times across dataset sizes. The testing dataset (50x50) runs almost instantaneously, whereas the small dataset (121x115) takes 12.137 ms. As matrix sizes increase, execution times rise considerably, with the medium dataset (550x480) requiring 984.676 ms, and the large dataset (962x1221) running for 11.166 seconds. The standard deviation values indicate varying levels of runtime consistency and with larger datasets, it shows greater fluctuation.

These results highlight the inefficiencies of the naïve implementation, especially for larger datasets where cache inefficiencies and memory access patterns significantly impact performance. Future optimizations should focus on improving parallelization and cache optimization.

8 References

1. GNU Make Manual: <https://www.gnu.org/software/make/manual/make.html>
2. IEEE POSIX Standard: <https://ieeexplore.ieee.org/servlet/opac?punumber=6880749>
3. YABMS Original Repository: <https://github.com/hawajkm/YABMS>
4. YABMS Forked Repository: <https://github.com/mohammed0x00/YABMS>