

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.

3.1 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 2: 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.

4 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.

4.1 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

4.2 Results

Table 1 summarizes the execution times for different dataset sizes.

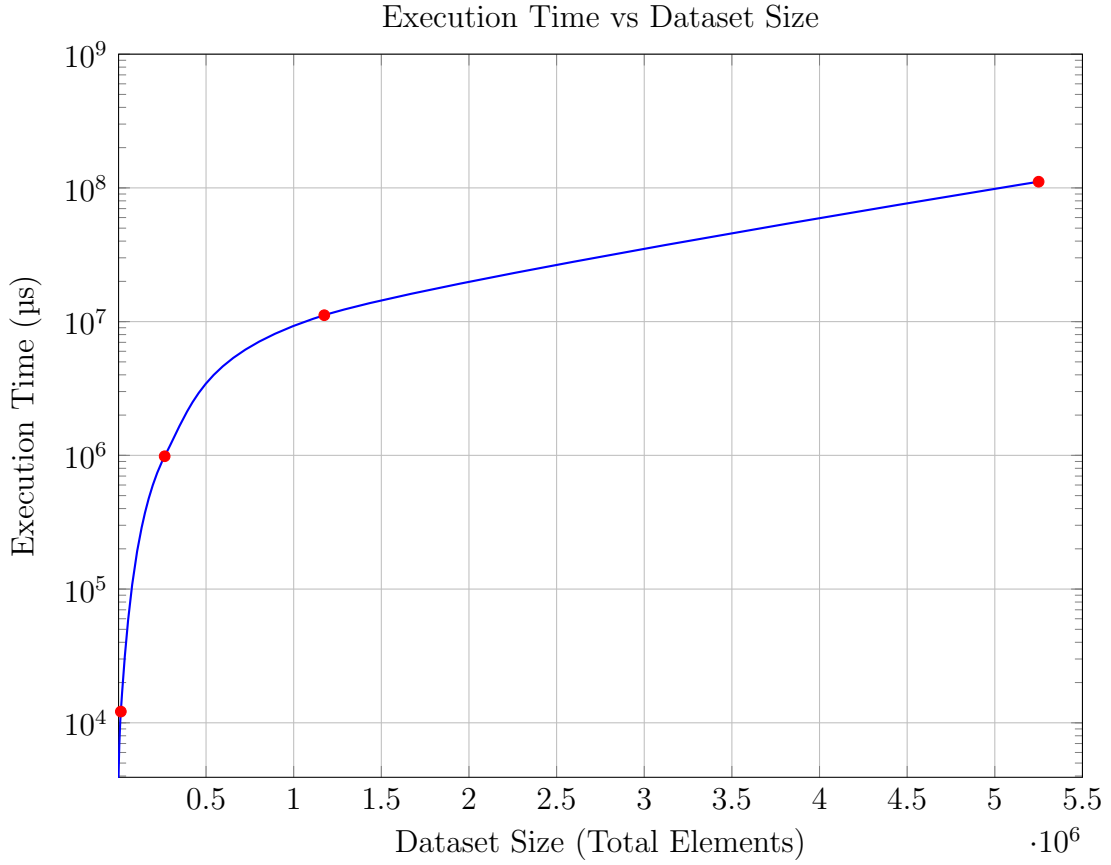
Dataset Size	Execution Time (ms)	Standard Deviation
Testing (16x8)	0.005	0.0008
Small (121x115)	12.137	1.189
Medium (550x480)	984.68	20.159
Large (962x1221)	11167	46.958
Native (2500x2100)	111219	661.744

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.

4.3 Graphical Representation

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



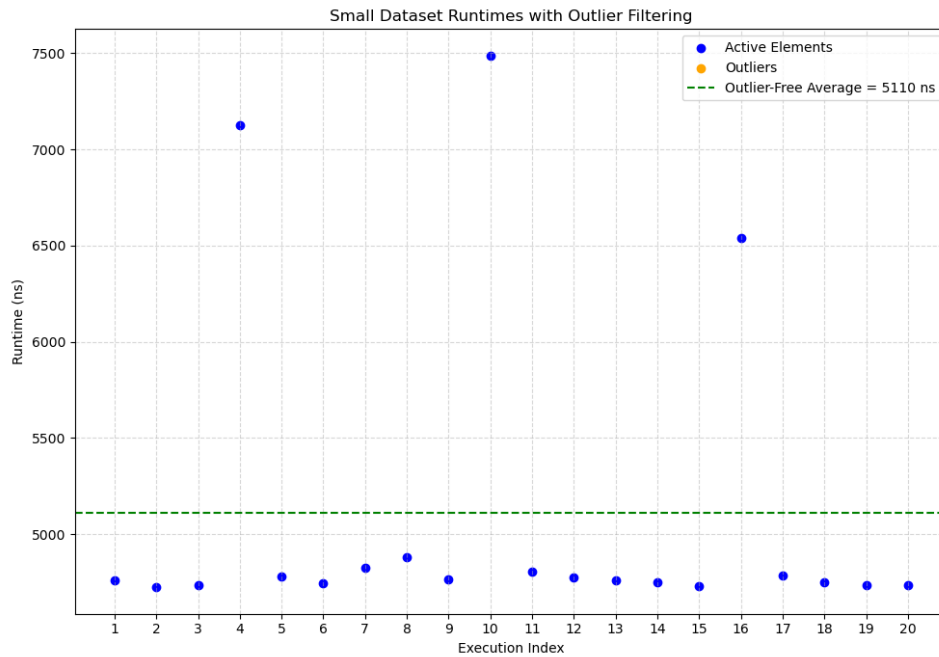
4.4 Detailed Results

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:

4.5 Testing Dataset Details

The following are the detailed results for the Testing dataset:

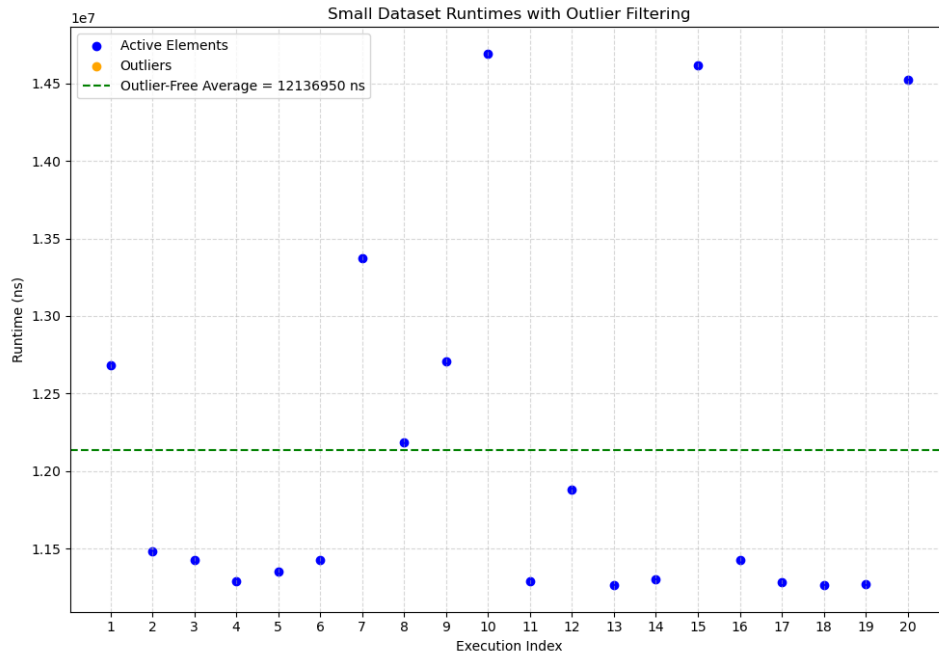


```
1 Setting up schedulers and affinity:
2   * Setting the niceness level:
3     -> trying niceness level = -20
4     + Process has niceness level = -20
5   * Setting up FIFO scheduling scheme and high priority ... Succeeded
6   * Setting up scheduling affinity ... Succeeded
7
8 Running "scalar_naive" implementation:
9   * Invoking the implementation 20 times .... Finished
10  * Verifying results .... Success
11  * Running statistics:
12    + Starting statistics run number #1:
13      - Standard deviation = 829
14      - Average = 5110
15      - Number of active elements = 20
16      - Number of masked-off = 0
17  * Runtimes (MATCHING): 5110 ns
18  * Dumping runtime informations:
19    - Filename: scalar_naive_runtimes.csv
20    - Opening file .... Succeeded
21    - Writing runtimes ... Finished
22    - Closing file handle .... Finished
```

Listing 3: Testing Dataset Output

4.6 Small Dataset Details

The following are the detailed results for the Small dataset:

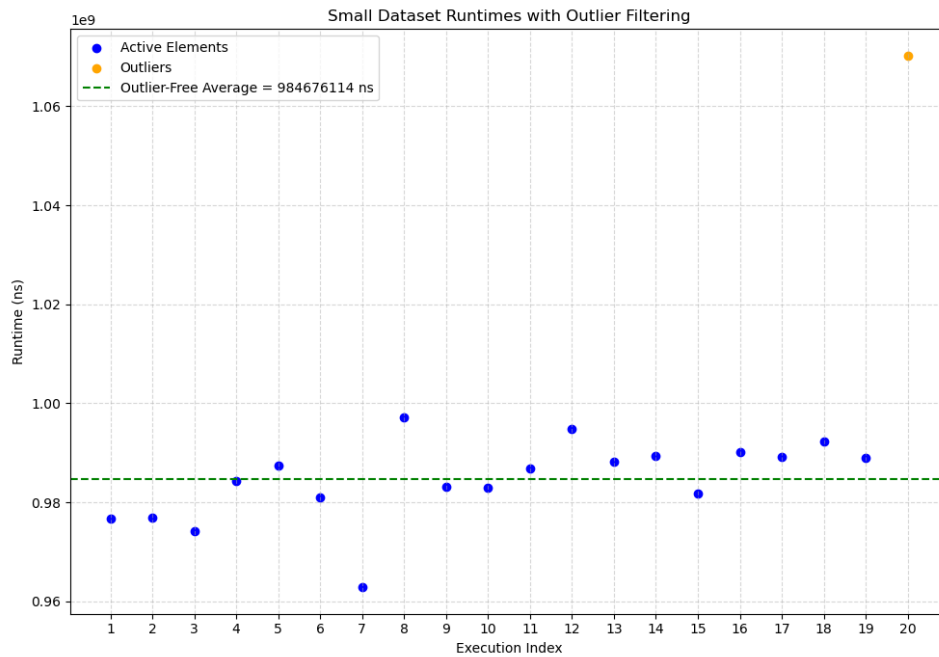


```
1 Setting up schedulers and affinity:
2   * Setting the niceness level:
3     -> trying niceness level = -20
4     + Process has niceness level = -20
5   * Setting up FIFO scheduling scheme and high priority ... Succeeded
6   * Setting up scheduling affinity ... Succeeded
7
8 Running "scalar_naive" implementation:
9   * Invoking the implementation 20 times .... Finished
10  * Verifying results .... Success
11  * Running statistics:
12    + Starting statistics run number #1:
13      - Standard deviation = 1189749
14      - Average = 12136950
15      - Number of active elements = 20
16      - Number of masked-off = 0
17  * Runtimes (MATCHING): 12136950 ns
18  * Dumping runtime informations:
19    - Filename: scalar_naive_runtimes.csv
20    - Opening file .... Succeeded
21    - Writing runtimes ... Finished
22    - Closing file handle .... Finished
```

Listing 4: Small Dataset Output

4.7 Medium Dataset Details

The following are the detailed results for the Medium dataset:

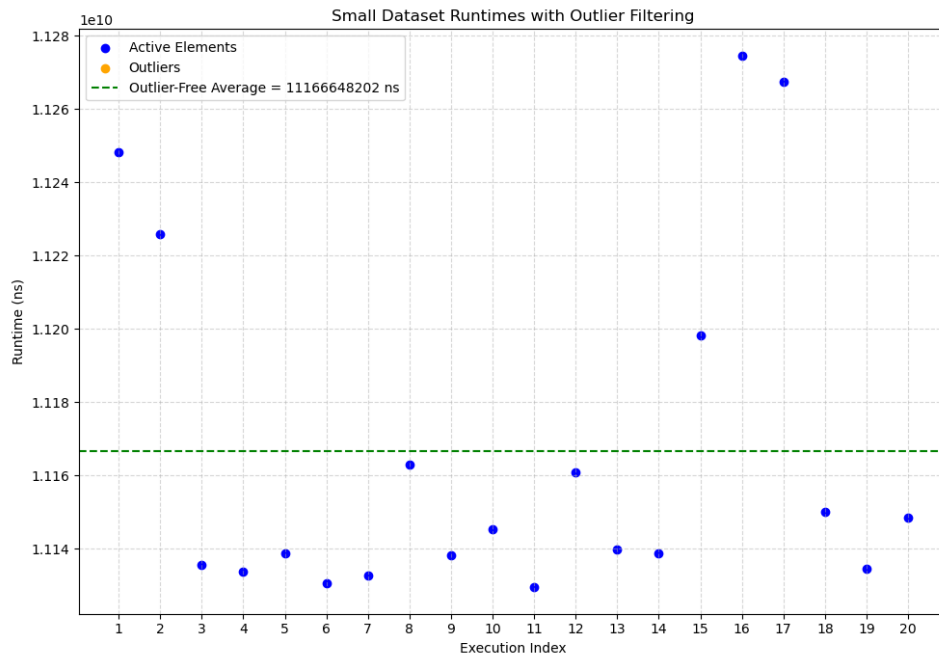


```
1 Setting up schedulers and affinity:
2   * Setting the niceness level:
3     -> trying niceness level = -20
4     + Process has niceness level = -20
5   * Setting up FIFO scheduling scheme and high priority ... Succeeded
6   * Setting up scheduling affinity ... Succeeded
7
8 Running "scalar_naive" implementation:
9   * Invoking the implementation 20 times .... Finished
10  * Verifying results .... Success
11  * Running statistics:
12    + Starting statistics run number #1:
13      - Standard deviation = 20159900
14      - Average = 988953986
15      - Number of active elements = 20
16      - Number of masked-off = 1
17    + Starting statistics run number #2:
18      - Standard deviation = 7861834
19      - Average = 984676114
20      - Number of active elements = 19
21      - Number of masked-off = 0
22  * Runtimes (MATCHING): 984676114 ns
23  * Dumping runtime informations:
24    - Filename: scalar_naive_runtimes.csv
25    - Opening file .... Succeeded
26    - Writing runtimes ... Finished
27    - Closing file handle .... Finished
```

Listing 5: Medium Dataset Output

4.8 Large Dataset Details

The following are the detailed results for the Large dataset:

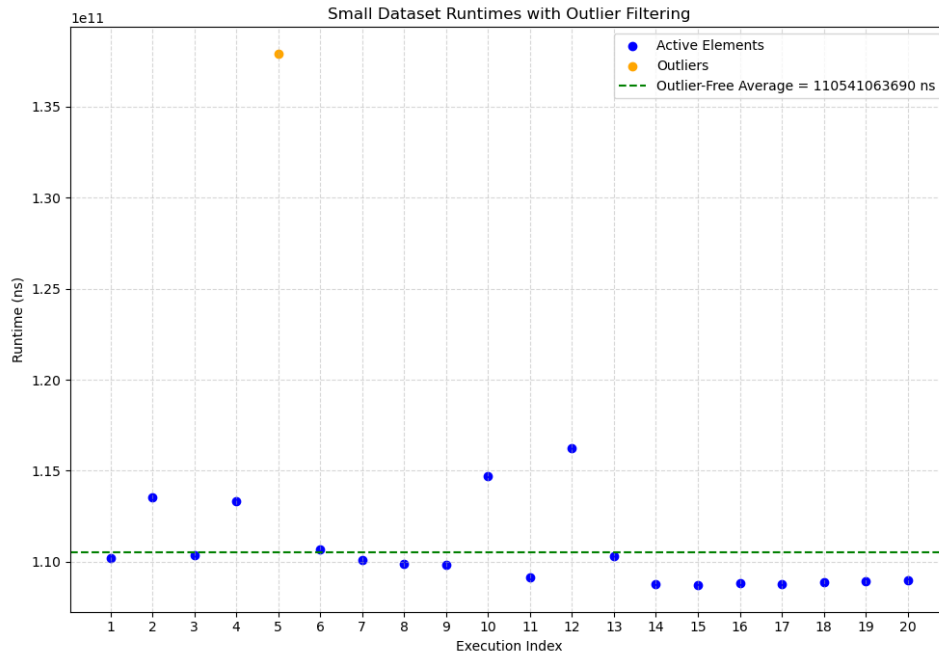


```
1 Setting up schedulers and affinity:
2   * Setting the niceness level:
3     -> trying niceness level = -20
4     + Process has niceness level = -20
5   * Setting up FIFO scheduling scheme and high priority ... Succeeded
6   * Setting up scheduling affinity ... Succeeded
7
8 Running "scalar_naive" implementation:
9   * Invoking the implementation 20 times .... Finished
10  * Verifying results .... Success
11  * Running statistics:
12    + Starting statistics run number #1:
13      - Standard deviation = 46958027
14      - Average = 11166648202
15      - Number of active elements = 20
16      - Number of masked-off = 0
17  * Runtimes (MATCHING): 11166648202 ns
18  * Dumping runtime informations:
19    - Filename: scalar_naive_runtimes.csv
20    - Opening file .... Succeeded
21    - Writing runtimes ... Finished
22    - Closing file handle .... Finished
```

Listing 6: Large Dataset Output

4.9 Native Dataset Details

The following are the detailed results for the Native 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_naive" implementation:
9   * Invoking the implementation 20 times ....
10
11
12 Finished
13   * Verifying results .... Success
14   * Running statistics:
15     + Starting statistics run number #1:
16       - Standard deviation = 661744362
17       - Average = 111909358290
18       - Number of active elements = 20
19       - Number of masked-off = 13
20     + Starting statistics run number #2:
21       - Standard deviation = 1411417123
22       - Average = 111218584463
23       - Number of active elements = 7
24       - Number of masked-off = 0
25   * Runtimes (MATCHING): 111218584463 ns
26   * Dumping runtime informations:
27     - Filename: scalar_naive_runtimes.csv
28     - Opening file .... Succeeded
29     - Writing runtimes ... Finished
30     - Closing file handle .... Finished
```

Listing 7: Native Dataset Output

5 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.

6 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>