# Memory Technologies

## INSTRUCTOR: PROF. HAMED FARBEH

## AMIRKABIR UNIVERSITY OF TECHNOLOGY
### (TEHRAN POLYTECHNIC)

Optimizing Algorithm Performance with HBM-PIM: A Matrix Multiplication
Case Study

Authors:

Morteza Adelkhani                                    Madelkhani@aut.ac.ir
Sara Zamani                                        sara.zamani73@aut.ac.ir

Student, project by:

Reza Adinepour                                       adinepour@aut.ac.ir

Spring 2024

# Evaluation, Academic Integrity and Submition

**Notes on the project:**

## 1- Evaluation

The evaluation of this project is based on the following:

a) **Algorithm Implementation (10%)**: Correctness and efficiency of the algorithm implementation.

b) **Simulation Setup (30%)**: Proper configuration and usage of the PIMSimulator.

c) **Performance Analysis (30%)**: Depth of analysis and understanding of performance metrics.

d) **Report Quality (20%)**: Clarity, completeness, and professionalism of the report.

e) **Presentation (10%)**: Effectiveness and clarity of the presentation.

## 2- Academic Integrity

Each student is expected to adhere to the highest standards of academic integrity. Any form of copying or plagiarism will result in severe penalties. Ensure your work is original and properly referenced.

## 3- Submission

The project submission guideline is as follows:

a) **Progress Reports**: Submit reports via the Courses portal until the deadline. Each delay in your submission is not acceptable.

b) **Final Report and Code**: Submit your final report and source code in a zipped folder.

c) **Presentation**: Present your findings in class during the final project meeting.

# Contents

# 1 Project Description

## 1.1 Proposed of this project:

Neurodegenerative diseases, including Alzheimer's In this project you are going to deal with processing in memory (PIM) structure. For surfing in real-world application of PIM, in this project you have to run an algorithm on real PIM device and report your result that you get. Based on the lack of accessibility of real PIM hardware you are going to using PIMSimulator which is based on HBM-PIM of Samsung.

## 1.2 Description of HBM-PIM:

High Bandwidth Memory (HBM) is a type of memory that's made to transfer data quickly and use less energy. It's built by stacking memory layers on top of each other, which lets them connect directly and move data fast. At the base of these layers, there's a special piece that controls the flow of data to and from other parts of the computer. HBM is often used with powerful computer chips like GPUs because it can handle a lot of data at once, making everything run smoother and faster.

Inside HBM, there are separate paths for data called pseudo-channels, and each one has smaller sections called banks where data is stored. When the computer needs to read data, it picks a specific path and bank, then grabs the data from there. This process is similar to how other types of memory work, but HBM's design lets it do this much quicker and with less energy.

Samsung has made a version of HBM called HBM-PIM that's even better because it can do some data processing right inside the memory itself. This means the computer doesn't have to move data around as much, which makes things faster and saves energy. This new design fits in with how memory is usually made, so it's easy to start using in products.

Overall, HBM and its improved version, HBM-PIM, are big steps forward for memory technology. They're really important for programs that need to process a lot of data quickly, like artificial intelligence and scientific computing, because they make everything more efficient and faster.

# 2 Project Detaile

As it described in previous section you will use the PIMSimulator to simulate the performance of different algorithms on the HBM-PIM architecture. Each student will be assigned one algorithm to analyze. The goal is to understand how PIM technology affects the performance and efficiency of these algorithms compared to traditional memory architectures. Your task to doing this project is as follow:

## 2.1 Choosing Algorithm:

You have to choose one of the following algorithms to analyze and submit it in your section of project table assignment in courses portal.

Note: The priority of choosing an algorithm is with the first student which choose it.

### 2.1.1   Sorting Algorithms:

a) QuickSort

b) MergeSort

c) HeapSort

d) Insertion Sort

e) Bubble Sort

### 2.1.2   Graph Algorithms:

a) Dijkstra's Algorithm

b) Breadth-First Search (BFS)

c) Depth-First Search (DFS)

d) Prim's Algorithm

e) Kruskal's Algorithm

### 2.1.3   Matrix Operations:

f) Matrix Multiplication

g) Sparse Matrix-Vector Multiplication (SpMV)

# IV. Machine Learning and Data Processing Algorithms:

h) K-means Clustering

## 2.2   Literature Review:

Each student will review existing research on HBM, PIM, and the specific algorithm assigned to them.

## 2.3   Algorithm Implementation:

Students will implement their assigned algorithm in a compatible programming language (e.g., C++, Python) if not already available.

## 2.4   Simulation Setup:

Students will set up the PIMSimulator to run their algorithms, configuring necessary parameters and optimizing the code to leverage PIM features.

## 2.5   Performance Analysis:

Students will run simulations to collect data on execution time, power consumption, and other relevant metrics.

## 2.6   Comparison:

Compare the performance of the algorithm on HBM-PIM with traditional memory architectures.

## 2.7   Report and Presentation:

Each student will compile their findings into a detailed report and present their results to the class.

The algorithm I have chosen for this project is matrix multiplication. Therefore, I first need to build the simulator, the installation steps of which I will explain below:

# 3   PIMSimulator

## 3.1   Overview

PIMSimulator is a cycle accurate model that Single Instruction, Multiple Data (SIMD) execution units that uses the bank-level parallelism in PIM Block to boost performance that would have otherwise used multiple times of bandwidth from simultaneous access of all bank. The simulator include memory and have embedded within it a PIM block, which consist of programmable command registers, general purpose register files and execution units.

Based on [github.com/umd-memsys/DRAMSim2](github.com/umd-memsys/DRAMSim2), the simulator includes the simulator includes

- PIM Block:

    - Register files including CRF (for command), GRF (for vector value), SRF (for scalar value)

    - ALU (ADD, MUL, MAC, MAD, MOVE, FILL, NOP, JUMP, EXIT)

- PIM Kernel:

    - Generate a set of memory transactions for enabling PIM operation

- HBM2 support (refer to `ini/HBM2_samsung_2M_16B_x64.ini`)

## 3.2   HW description

PIM is a HBM stack that is pin compatible with HBM2 and have embedded within it a PIM block

---

Figure 1: Architecture of PIM

### 3.2.1 Base Architecture

- Each channel is logically independent memory, so it has a dedicated independent controller.

- Read [Addr], Write [Addr]

- Activate, Read, Write, Precharge, Refresh, `Activate_pim`, `ALU_pim`, `Precharge_pim`, `READ_pim`

1. HBM2

   (a) System Specification: `system_hbm.ini`

   (b) HBM Specification: `ini/HBM2_samsung_2M_16B_x64.ini`

      i. 1 PIM block per 2 banks, 4 Bank per Bankgroup, 4 Bank group per pseudo channel, 4 pseudo channel per die, 4 die per stack.

      ii. Prefetch size : 256bit

      iii. burst length: 4n

      iv. Pin speed: 2Gbps

      v. The simulator supports the pseudo-channel mode only, and we assume that each pseudo-channel is totally independent.

### 3.2.2 Address mapping

1. The address mapping is used when the memory controller decodes the address from host.

2. Use Scheme8 addressing mode for PIM functionality.

   ```
   |<-rank->|<-row->|<-col high->|<-bg->|<-bank->|<-chan->|<-col low->|<-offset
   ->|
   ```

3. the length of `col_low` is log(BL * JEDEC_DATA_BUS_BUTS/8), which are 5b both for HBM2

4. You can also change the current addressing mode dynamically (Not recommended, though)

   ```
   // Static Setting in system_*.ini
   ADDRESS_MAPPING_SCHEME=Scheme8
   ```

Figure 2: HBM2 case

### 3.2.3   PIM Block Placement

`BANKS_PER_PIM_BLOCK = NUM_BANKS / NUM_PIM_BLOCKS`

1. if `2 * NUM_PIM_BLOCKS == NUM_BANK`, a PIM block is located per two banks.

   (a) `NUM_PIM_BLOCKS = 8, NUM_BANKS = 16`

2. if `NUM_PIM_BLOCKS == NUM_BANKS`, a PIM Block (PB) is located per banks.

   (a) `NUM_PIM_BLOCKS = 8, NUM_BANKS = 8`

- Supports RISC-style 32-bit instructions

- Three instructions types

   - 4 Arithmetic: ADD, MUL, MAC, MAD
   - 2 Data transfer: MOV, FILL
   - 3 control flows: NOP, JUMP, EXIT

- JUMP instruction

   - Zero-cycle static branch: supports only a pre-programmed numbers of iterations

- Operand type:

   - Vector Register (GRF_A, GRF_B)
   - Scalar Register (SRF)
   - Bank Row Buffer

- PIM instructions are stored in the Command Register File (CRF), and memory command triggers a CRF to perform a target instruction

| Type | Command | Description | Result (DST) | Operand (SRC0) | Operand (SRC1) |
|------|---------|-------------|--------------|----------------|----------------|
| Arithmetic | ADD | addition | GRF | GRF, BANK, SRF | GRF, BANK, SRF |
| Arithmetic | MUL | multiplication | GRF | GRF, BANK, SRF | GRF, BANK, SRF |
| Arithmetic | MAC | multiply-accumulate | GRF_B | GRF, BANK | GRF, BANK, SRF |
| Arithmetic | MAD | multiply-and-add | GRF | GRF, BANK | GRF, BANK, SRF |
| Data | MOV | load or store data from register to bank | GRF, SRF | GRF, SRF | |
| Data | FILL | copy data from bank to register | GRF, BANK | GRF, BANK | |
| Control | NOP | do nothing | | | |
| Control | JUMP | jump instruction | | | |
| Control | EXIT | exit instruction | | | |

Table 1: Table of Commands

- − each memory command increments the CRF PC

- DRAM commands decide where to retrieve data from DRAM for PIM arithmetic operations

### 3.2.4 Movement of Data

| Mode | Transaction | PIM Instruction | Operation |
|------|-------------|-----------------|-----------|
| SB | Read | - | Normal Memory Read |
| SB | Write | - | Normal Memory Write |
| HAB | Write | - | PIM Write (Host to PIM Register) |
| PIM | - | MOV | read or write from bank to PIM Register |
| PIM | - | FILL | write from bank to PIM Registers |

Table 2: Table of Transactions and Operations

1. SB mode: standard DRAM operation

2. HAB mode: Allowing concurrent accesses to multiple banks with a single DRAM command

3. PIM mode: Triggers the execution of PIM commands on the CRF by DRAM Command

## 3.3   Setup

## 3.4   Prerequisites

- `Scons` tool for compiling PIMSimulator:

  ```
  $ sudo apt install scons
  ```

- `gtest` for running test cases:

  ```
  $ sudo apt install libgtest-dev
  ```

## 3.5   Installing

- To Install PIMSimulator:

  ```
  # compile
  scons
  ```

After entering the `scons` command, we encountered the following errors.



Figure 3: Error in build

After some investigation, I realized that the problem was with the compiler. I had both `gcc` and `g++` compilers installed on my system, and by default, the `gcc` compiler was selected. I changed the default compiler to `g++` with the following command, and the error problem was resolved.

```
$ sudo update-alternatives --install /usr/bin/g++ g++ /usr/bin/g++-9 60
$ sudo update-alternatives --config g++
```

With the compiler change, the build is completed successfully



Figure 4: Done building targets

In this simulator, there are pre-written examples that we can list as follows:

## 3.6   Launch a Test Run

- Show a list of test cases

```
./sim --gtest_list_tests

# Example
PIMKernelFixture.
gemv_tree
gemv
mul
add
relu
MemBandwidthFixture.
hbm_read_bandwidth
hbm_write_bandwidth
```

```
     PIMBenchFixture.
     gemv
     mul
     add
     relu
```

For further investigation, we will thoroughly examine and explain the matrix multiplication code file located in the path /PIMSimulator/data/gemv/gen_gemv.py.

# 4 Matrix-vector Multiplication Code Breakdown

## 4.1 Introduction

This code explains a Python script that performs matrix-vector multiplication. The script involves generating random input data and weights, performing matrix multiplications, and saving the results to files.

```python
import numpy as np

# min dim_in = 128 -> 256bit / 16bit
# min dim_out = 8 PIM block
BATCH = 1
REAL_DIM_IN = 1024
DIM_IN = 1024
DIM_OUT = 4096

np.set_printoptions(precision=20)
np.random.seed(41)

batch_in = np.random.standard_normal(size=(DIM_IN, BATCH)).astype('float16')
for i in range(REAL_DIM_IN, DIM_IN):
  for j in range(0, BATCH):
    batch_in[i][j] = 0

data_w = np.random.standard_normal(size=(DIM_OUT, DIM_IN)).astype('float16')

np.random.shuffle(data_w)
batch_out = np.zeros((DIM_OUT, BATCH)).astype('float16')
batch_out = np.matmul(data_w, batch_in)

batch_out2 = np.zeros((DIM_OUT, BATCH)).astype('float16')

for y in range(0, DIM_OUT):
  for x in range(0, DIM_IN):
    batch_out2[y] += data_w[y][x] * batch_in[x][0]

batch_in = batch_in.T.copy()
batch_out = batch_out.T.copy()
batch_out2 = batch_out2.T.copy()

np.save("gemv_input_" + str(DIM_OUT) + "x" + str(DIM_IN), batch_in)
np.save("gemv_weight_" + str(DIM_OUT) + "x" + str(DIM_IN), data_w)
```

```
36   np.save("gemv_output_" + str(DIM_OUT) + "x" + str(DIM_IN), batch_out)
37   np.save("test_output_" + str(DIM_OUT) + "x" + str(DIM_IN), batch_out2)
38   print(batch_in)
39   print(batch_out)
40   print(batch_out2)
41   print(batch_in.shape)
42   print(batch_out.shape)
```

Listing 1: Python Code

## 4.2   Set constants

```
1   BATCH = 1
2   REAL_DIM_IN = 1024
3   DIM_IN = 1024
4   DIM_OUT = 4096
```

- BATCH: The batch size of the input vectors.

- REAL_DIM_IN: Actual input dimension size.

- DIM_IN: Padded input dimension size (could be the same or larger than REAL_DIM_IN).

- DIM_OUT: Output dimension size.

## 4.3   Generate input data

```
1   batch_in = np.random.standard_normal(size=(DIM_IN, BATCH)).astype('float16')
```

Create a DIM_IN x BATCH matrix filled with random values from a standard normal distribution, converted to float16 precision.

## 4.4   Zero out the padded part of the input

```
1   for i in range(REAL_DIM_IN, DIM_IN):
2     for j in range(0, BATCH):
3       batch_in[i][j] = 0
```

Set the elements of batch_in to zero for indices from REAL_DIM_IN to DIM_IN to handle the padded part.

## 4.5   Generate weight data

```
1   data_w = np.random.standard_normal(size=(DIM_OUT, DIM_IN)).astype('float16')
2   np.random.shuffle(data_w)
```

Create a DIM_OUT x DIM_IN matrix filled with random values from a standard normal distribution, converted to float16 precision. Shuffle the rows of data_w.

## 4.6   Initialize output matrices and perform matrix multiplication

```
1   batch_out = np.zeros((DIM_OUT, BATCH)).astype('float16')
2   batch_out = np.matmul(data_w, batch_in)
```

## 4.7   Manual matrix multiplication for verification

```
1   batch_out2 = np.zeros((DIM_OUT, BATCH)).astype('float16')
2
3   for y in range(0, DIM_OUT):
4     for x in range(0, DIM_IN):
5       batch_out2[y] += data_w[y][x] * batch_in[x][0]
```

`batch_out2` is calculated element-wise by iterating over each element of the result and summing the product of corresponding elements from `data_w` and `batch_in`.

## 4.8   Transpose the input and output matrices

```
1   batch_in = batch_in.T.copy()
2   batch_out = batch_out.T.copy()
3   batch_out2 = batch_out2.T.copy()
```

## 4.9   Save the input, weight, and output matrices to files

```
1   np.save("gemv_input_" + str(DIM_OUT) + "x" + str(DIM_IN), batch_in)
2   np.save("gemv_weight_" + str(DIM_OUT) + "x" + str(DIM_IN), data_w)
3   np.save("gemv_output_" + str(DIM_OUT) + "x" + str(DIM_IN), batch_out)
4   np.save("test_output_" + str(DIM_OUT) + "x" + str(DIM_IN), batch_out2)
```

## 4.10   Print the matrices and their shapes

```
1   print(batch_in)
2   print(batch_out)
3   print(batch_out2)
4   print(batch_in.shape)
5   print(batch_out.shape)
```

After running the code, the output is as follows, and the matrices are saved as Numpy arrays in the same directory.

Now, after preparing the matrices, it is time to perform the matrix multiplication once with PIM and once without PIM.

To do this, we enter the following command in the simulator directory. The simulator runs the program once without considering PIM and once with considering PIM, and displays the outputs.

```
1   $ ./sim --gtest_filter=PIMBenchFixture.gemv
```
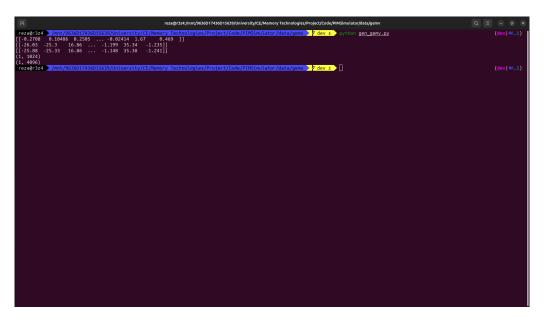
Figure 5: Output of `gen_gemv.py`

# 5   Output simulation
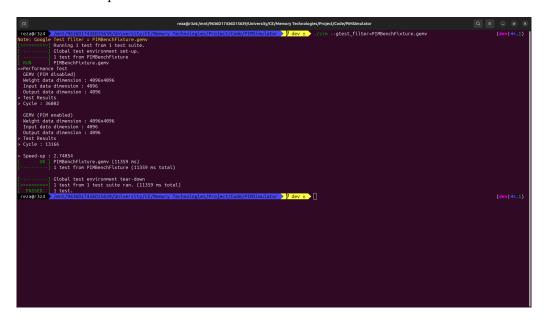
The simulation output is as follows:



Figure 6: Output of PIM simulation

## 5.1   Analysis

### 5.1.1   Test Context

The test involves one test suite named `PIMBenchFixture` with one test case `gemv`. The simulation was executed using Google Test with the filter `PIMBenchFixture.gemv`.

### 5.1.2 Performance Test

The GEMV operation was performed twice:

1. **Without PIM:**

   - **Weight Data Dimension:** 4096x4096
   - **Input Data Dimension:** 4096
   - **Output Data Dimension:** 4096
   - **Cycles:** 36,082

2. **With PIM:**

   - **Weight Data Dimension:** 4096x4096
   - **Input Data Dimension:** 4096
   - **Output Data Dimension:** 4096
   - **Cycles:** 13,166

### 5.1.3 Results Interpretation

- **Cycle Count:**

  - **Without PIM:** The matrix multiplication took 36,082 cycles.
  - **With PIM:** The matrix multiplication took 13,166 cycles.

- **Speed-up:** The speed-up achieved by enabling PIM is calculated as 2.74054.

### 5.1.4 Overall Test Duration

The total time taken for the test execution was 11,359 milliseconds.

### 5.1.5 Significant Performance Improvement

The use of PIM technology resulted in a substantial reduction in the number of cycles required to complete the GEMV operation. Specifically, PIM enabled a 2.74x speed-up compared to the traditional approach without PIM. This indicates that PIM technology can significantly enhance the performance of matrix-vector multiplication by reducing the cycle count, which directly correlates to faster computation times.

### 5.1.6 Consistency in Data Dimensions

The dimensions of the weight, input, and output data were consistent across both tests (4096x4096 for weight and 4096 for both input and output). This ensures a fair comparison between the PIM-enabled and PIM-disabled scenarios.

### 5.1.7 Total Execution Time

The entire test suite, including setup and teardown, completed in approximately 11.359 seconds. This time includes not only the GEMV operations but also any additional overhead associated with the test framework and simulation environment.

### 5.1.8 Pass/Fail Status

The test passed successfully, indicating that the GEMV functionality works as expected in both PIM-enabled and PIM-disabled modes.

# 6  Conclusion

The simulation results clearly demonstrate the performance benefits of utilizing PIM technology for matrix-vector multiplication tasks. By significantly reducing the number of cycles required for the GEMV operation, PIM can lead to faster computations and more efficient processing, making it a valuable approach for high-performance computing applications.

The successful completion of this test and the observed speed-up highlight the potential of PIM technology in accelerating matrix operations, which are fundamental in various computational workloads, including machine learning and scientific computing.

# References

[1] SAITPublic. Pimsimulator. https://github.com/SAITPublic/PIMSimulator. Accessed: 2024-07-20.

[2] Samsung Advanced Institute of Technology. Samsung advanced institute of technology. https://www.sait.samsung.co.kr/. Accessed: 2024-07-20.

[3] UMD Memsys. Dramsim2. https://github.com/umd-memsys/DRAMSim2. Accessed: 2024-07-20.