





# Exploring the Potential of Real PIM Architecture for Quantum Circuit Simulation

Dongin Lee\* (National University of Singapore), Enhyeok Jang\*, Seungwoo Choi, Junwoong An, Cheolhwan Kim, and Won Woo Ro (Yonsei University)

\*These two authors contributed equally to this work

E-mail: dongin.lee@u.nus.edu, {enhyeok.jang, wro}@yonsei.ac.kr

Session 3D-3 (T3-2) Frameworks and Modeling for Computing-In-Memory, Room Mars/Mercury

Tuesday, January 21st, 2025

#### Contents

- Background and Motivation
- **PIMUTATION Design**
- Experimental Methodology
- Result and Analysis
- Conclusion



#### Contents

- Background and Motivation
- PIMUTATION Design
- Experimental Methodology
- Result and Analysis
- Conclusion



## Why Quantum Circuit Simulation is important ?

- Current quantum computers have low accuracy, long job wait times, and high execution fee.
  ✓ QCS is advantageous to <u>study computational processes of quantum</u> processors alernatively.
- QCSs are also essential for the <u>verification of intermediate states</u> in quantum algorithms,
  ✓ Real QCs lose their superposed/entangled states after qubit measurements.



- State vector-based QCS is the most common method.
- State vectors for n-qubit system requires memory capacity of  $2^n$  elements of complex amplitudes.
  - ✓ EX) <u>1Q gate for 10-qubit system</u> correspond to 512 times of <u>2x2 Matrix 2x1 Vector Multiplications</u>.



- State vector-based QCS is the most common method.
- State vectors for n-qubit system requires memory capacity of  $2^n$  elements of complex amplitudes.
  - ✓ EX) <u>1Q gate for 10-qubit system</u> correspond to 512 times of <u>2x2 Matrix 2x1 Vector Multiplications</u>.



- State vector-based QCS is the most common method.
- State vectors for n-qubit system requires memory capacity of  $2^n$  elements of complex amplitudes.
  - ✓ EX) <u>1Q gate for 10-qubit system</u> correspond to 512 times of <u>2x2 Matrix 2x1 Vector Multiplications</u>.



000
001
010
011
100
101
110
111

- State vector-based QCS is the most common method.
- State vectors for n-qubit system requires memory capacity of  $2^n$  elements of complex amplitudes.
  - ✓ EX) <u>1Q gate for 10-qubit system</u> correspond to 512 times of <u>2x2 Matrix 2x1 Vector Multiplications</u>.



000
001
010
011
100
101
110
111

#### Trends in Quantum Circuit Simulation in Industry Fields



• IBM, Google, and NVIDIA have competitively reported large-scale quantum circuit simulations leveraging multiple GPU systems (V100 or A100).

9

#### Trends in Quantum Circuit Simulation in Industry Fields



https://www.youtube.com/watch?v=T-a\_rCfKTE https://quantumai.google/qsim/choose\_hw https://www.youtube.com/watch?v=Qdb3xrzT1gc

#### Workload Characteristics of QCSs

- QCS consists of iterations of small-sized matrix-vector multiplications.
- Each gate requires *distinct access patterns* to all  $2^n$  amplitudes at least once.
  - ✓ Low locality data access, which Incurs a lot of cache misses.
  - ✓ Frequent data movements between CPU and main memory.
- Various quantum applications exhibit memory (bandwidth)-bound properties.





#### Workload Characteristics of QCSs

- QCS consists of iterations of small-sized matrix-vector multiplications.
- Each gate requires *distinct access patterns* to all  $2^n$  amplitudes at least once.
  - ✓ Low locality data access, which Incurs a lot of cache misses.
  - ✓ Frequent data movements between CPU and main memory.
- Various quantum applications exhibit memory (bandwidth)-bound properties.



< Roofline Analysis on Intel Xeon Silver 4215 Processor >

#### Workload Characteristics of QCSs

- QCS consists of iterations of small-sized matrix-vector multiplications.
- Each gate requires *distinct access patterns* to all  $2^n$  amplitudes at least once.
  - ✓ <u>Low locality</u> data access, which Incurs a lot of cache misses.
  - ✓ Frequent data movements between CPU and main memory.
- Various quantum applications exhibit memory (bandwidth)-bound properties.



< Roofline Analysis on Intel Xeon Silver 4215 Processor >

#### Design Challenges to Run QCS with UPMEM

- QCSs require FP Ops for calculating complex-valued amplitudes.
  - ✓ However, DPU (UPMEM's processing unit) only <u>supports integer Ops</u>.
  - ✓ FP Ops are implemented using software libraries, but, which <u>degrades the performance significantly</u>. ☺
- Depending on the gate's positions, the memory access patterns can vary greatly.
  ✓ However, there is no direct communication channel between DPUs. ⊗



#### Contents

- Background and Motivation
- **PIMUTATION Design**
- Experimental Methodology
- Result and Analysis
- Conclusion



#### **PIMUTATION:** <u>PIM</u> for <u>qUanTum</u> circuit Simul<u>ATION</u>



First, the host CPU sends quantum *gate* information and initialized *state vectors* to the DPU.

To minimize *FP Ops and multiplication* operations, PIMutation utilizes Gate Merging and Row Swapping.

To reduce *communication between DPUs,* it is divided into sub-circuits and transmitted to each DPU if the original QC is separable.

Finally, sub-circuit results are reconfigured in the host CPU and then restored to the full state vector, suppressing communication between DPUs.

#### **PIMUTATION Optimization Techniques**

- Gate merging: Exploiting UPMEM native Integer Ops
  - $\checkmark$  This replaces irrational number divisions in gates with bit-wise shifts that UPMEM can process.
- State Resizing: Maintaining data type of amplitudes as Integer
  ✓ By normalizing the min. amplitude to 1, the amplitudes can be integers during the entire process.
- Row swapping: Avoiding matrix multiplication for CX, SWAP gates
  - ✓ For matrices with only 0 or 1 elements, amplitudes are processed by move Ops, instead of MUL.
- Vector Partitioning: Partitioning separable quantum programs, if possible
  - $\checkmark$  Sub-circuits are isolated each other and processed at each DPU w/o communication between DPUs.



#### Gate Merging

• Quantum gates often require irrational division.

✓ For example,  $RX(\frac{\pi}{2}) = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -i \\ -i & 1 \end{bmatrix}$ ,  $RY(\frac{\pi}{2}) = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$ , and  $H = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ 

- GM provides an avoidance of FP Ops.
- Merged gates can be processed just using integer Ops or bit-wise Ops.



< Bernstein-Vazirani Circuit for the secret string of 111 >

$$\mathsf{H} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

#### Gate Merging

• Quantum gates often require irrational division.

✓ For example,  $RX(\frac{\pi}{2}) = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -i \\ -i & 1 \end{bmatrix}$ ,  $RY(\frac{\pi}{2}) = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$ , and  $H = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ 

- GM provides an avoidance of FP Ops.
- Merged gates can be processed just using integer Ops or bit-wise Ops.



## **Row Swapping**

- Some gates (CX or SWAP) are just switching amplitudes.
- These are implemented by directly swapping the corresponding row pairs.
  ✓ FP Mult Ops (hard to UPMEM ☺) => Move arrays directly



## **Vector Partitioning**

- Separable quantum circuit can be independently processed into sub-circuits.
- They are handled at each DPU and reconstructed to original vector on CPU.



#### Contents

- Background and Motivation
- PIMUTATION Design
- Experimental Methodology
- Result and Analysis
- Conclusion



#### Backend: Real UPMEM PIM Server (provided by UPMEM Co.)

- A single server provides 20 UPMEM DIMMs
- Each UPMEM DIMM includes 8 memory chips per rank
- Each memory chip is coupled with 8 processing element
- Processing element is called DRAM Processing Unit (DPU)



#### **Experimental Set-Up**

#### **Table 1: The Platform for Experiments**

#### **Table 2: Evaluated Benchmark Quantum circuits**

CPU Model	Intel(R) Xeon(R) Silver 4215
Sockets	2
Mamory	4 x 64 GB DDR4-2666
Wiemory	RDIM Dual Rank DRAM (256 GB)
PIM Memory	$20 \times DDR4-2400$ PIM Modules (160 GB)
<b>UPMEM SDK version</b>	2023.2.0
QuEST version	v3.7.0 [15]
Compiler version	gcc 8.3.0
CMake version	3.13.4
Energy Measurement	Intel RAPL [7]
	55

Algorithm [Reference]	Qubits	# of 1Q	# of 2Q	Abbr.
		Gates	Gales	
BB84 Protocol [5]	n	2n	0	BB_n
Bernstein–Vazirani [6]	n	2n	n-1	BV_n
Error Detection Code [9]	n	2n	2n-2	EDC_n
Hidden Subgroup	22	6n	2n	HS_2n
Problems [20]				
Quantum Random		n n	0	QRNG_n
Number Generator [20]				
Exclusive-OR [25]	n	0	n-1	XOR_n

Based on Float64, UPMEM PIM memory can <u>store up to 33-qubit systems</u>' full

state-vector, which requires 128 GB.

## **Evaluated Scenarios**

- Single-DPU experiments
  - Evaluating the performance for:
    - ✓ Gate Merging (GM) + State Resizing Included
    - ✓ Row Swapping (RS)
- Multi-DPU experiments
  - Evaluating the performance for:
    - ✓ Gate Merging (GM ) + State Resizing Included
    - ✓ Row Swapping (RS)
    - ✓ Vector Partitioning (VP)
  - Evaluating Speed-up and Energy consumption by PIMutation over the *C language-based cutting-edge QuEST simulator* on CPU
  - Evaluations on large-scale (16- and 32-qubit) benchmarks

#### Contents

- Background and Motivation
- PIMUTATION Design
- Experimental Methodology
- Result and Analysis
- Conclusion



#### **Single-DPU Performance**



- In 2-qubit benchmarks (XX\_2), Gate Merging (GM) and Row Swapping (RS) achieves a speedup of 0.26x and 0.13x over Baseline (No Optimization).
- In 4-qubit benchmarks (XX\_4), Gate Merging (GM) and Row Swapping (RS) achieves a speedup of 2.99x and 1.65x over Baseline (No Optimization).

## **Multi-DPU Performance**



- 16-qubit Workloads
  - ✓ 3.42x and 2.56x speedup over the QuEST in 4 and 8 DPUs, respectively.
  - $\checkmark$  25% energy saving over the QuEST in average of 4 and 8 DPUs.
- 32-qubit Workloads
  - ✓ 16.64x and 16.37x speedup over the QuEST in 8 and 16 DPUs, respectively.
  - $\checkmark$  75% energy saving over the QuEST in average of 8 and 16 DPUs.

## **Multi-DPU Performance**



- More DPUs do not always provide better performance.
- More DPUs can leverage parallelism, but its increased reconfiguration overhead can degrade overall performance.

 ✓ For example, 4 DPUs perform better than 8 DPUs.

=> Estimating the optimal number of DPUs for QCS would be an interesting future work.

- 16-qubit Workloads
  - ✓ 3.42x and 2.56x speedup over the QuEST in 4 and 8 DPUs, respectively.
  - $\checkmark$  25% energy saving over the QuEST in average of 4 and 8 DPUs.
- 32-qubit Workloads
  - ✓ 16.64x and 16.37x speedup over the QuEST in 8 and 16 DPUs, respectively.
  - $\checkmark$  75% energy saving over the QuEST in average of 8 and 16 DPUs.

#### Contents

- Background and Motivation
- PIMUTATION Design
- Experimental Methodology
- Result and Analysis
- Conclusion



## Conclusion

- We note that QCSs are <u>bandwidth-bound</u> workloads.
  - ✓ If we adopt PIM, therefore, *fast and low-energy* quantum circuit simulations would be achieved.
- We leverage UPMEM, the first commercialized PIM, to verify our observation.
  ✓ UPMEM can only support <u>Integer Ops natively</u> and <u>communication between DPUs is not available</u>.
  ✓ These design challenges should be considered when implementing quantum circuit simulations.
- PIMUTATION can minimize FP Ops and communications between DPUs. We propose 4 optimization techniques.
  - Gate Merging
  - State Resizing
  - Row Swapping
  - Vector Partitioning
- PIMUTATION achieves <u>faster and lower energy consumption</u>, compared to the state-of-the-art high-speed QuEST simulator.

#### **PIMUTATION Contributors**







#### **Thank You Your Time!**

Q & A

