BMQSim: Overcoming Memory Constraints in Quantum Circuit Simulation with a High-Fidelity Compression Framework

Boyuan Zhang

Indiana University Bloomington, USA bozhan@iu.edu

Bo Fang*

Pacific Northwest National Laboratory Richland, USA bo.fang@pnnl.gov

Fanjiang Ye

Indiana University Bloomington, USA fanjye@iu.edu

Luanzheng Guo

Pacific Northwest National Laboratory Richland, USA lenny.guo@pnnl.gov

Fengguang Song

Indiana University Bloomington, USA fgsong@iu.edu Nathan Tallent Pacific Northwest National Laboratory Richland, USA Nathan.Tallent@pnnl.gov

Dingwen Tao*

Indiana University Bloomington, USA taodingwen@ict.ac.cn

Abstract

Full-state quantum circuit simulation requires exponentially increased memory size to store the state vector as the number of qubits scales, presenting significant limitations in classical computing systems. Our paper introduces BMQSIM, a novel state vector quantum simulation framework that employs lossy compression to address the memory constraints on graphics processing unit (GPU) machines. BMQSIM effectively tackles four major challenges for state-vector simulation with compression: frequent compression/decompression, high memory movement overhead, lack of dedicated error control, and unpredictable memory space requirements. Our work proposes an innovative strategy of circuit partitioning to significantly reduce the frequency of compression occurrences. We introduce a pipeline that seamlessly integrates compression with data movement while concealing its overhead. Additionally, BMQSIM incorporates the first GPU-based lossy compression technique with point-wise error control. Furthermore, BMQSIM features a two-level memory management system, ensuring efficient and stable execution. Our evaluations demonstrate that BMQSIM can simulate the same circuit with over 10 times less memory usage on average, achieving fidelity over

This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License. *ICS '25, Salt Lake City, UT, USA* © 2025 Copyright held by the owner/author(s). ACM ISBN 979-8-4007-1537-2/25/06 https://doi.org/10.1145/3721145.3725747 0.99 and maintaining comparable simulation time to other state-of-the-art simulators.

CCS Concepts

 \bullet Theory of computation \rightarrow Data compression; \bullet Computing methodologies \rightarrow Quantum mechanic simulation.

Keywords

Quantum simulation, GPU, lossy compression, memory footprint.

ACM Reference Format:

Boyuan Zhang, Bo Fang, Fanjiang Ye, Luanzheng Guo, Fengguang Song, Nathan Tallent, and Dingwen Tao. 2025. BMQSim: Overcoming Memory Constraints in Quantum Circuit Simulation with a High-Fidelity Compression Framework. In 2025 International Conference on Supercomputing (ICS '25), June 08–11, 2025, Salt Lake City, UT, USA. ACM, New York, NY, USA, 16 pages. https: //doi.org/10.1145/3721145.3725747

1 Introduction

Quantum computing has emerged as a significant paradigm within the High Performance Computing (HPC) community. Its unique characteristics have drawn considerable attention. In recent years, quantum computing has proven effective in addressing key problems across various fields, such as machine learning [4, 36, 47], quantum chemistry [1], optimization problems [13], and financial modeling [45]. The advancement of quantum hardware aligns with the increasing impact of quantum computing. For instance, the state-of-the-art (SOTA) IBM Condor quantum system now

^{*}Co-corresponding authors.

supports 1,121 qubits, more than double of the 433 qubits supported by last year's Osprey quantum system [15]. Executing quantum algorithms on real quantum computers, however, faces fundamental challenges. First, in the current Noisy Intermediate-Scale Quantum (NISQ) era [44], noise interference in the hardware results in inaccurate measurement distribution. Second, designing new quantum algorithms requires iterative trials to verify, which is impractical on quantum computer platforms. Third, publicly available quantum computers (specifically those with a large number of qubits, e.g., > 16) are much less resourceful and usually reside in cloud services; hence, access to those machines is limited. Thus, quantum circuit simulation has become an essential approach for realizing the full potential of quantum computing [25]. Running a full-state quantum circuit simulation (i.e., state-vector simulation) presents a formidable challenge: as the number of simulated gubits increases, the memory requirement grows exponentially. Several significant issues are associated with this: (1) Simulating large quantum systems requires extensive memory capacity in classical systems. For instance, simulating a 48-qubit circuit would fully occupy the entire memory of Frontier (4.6 petabytes of DDR4 memory), the most advanced HPC machine currently available [2]. (2) Even when the memory capacity requirement is met, accessing such HPC systems requires dedicated allocation, which is usually quite competitive due to high demand. Consequently, researchers in quantum computing are often constrained to work with much smaller machines, such as personal computers or local workstations that typically have only dozens of gigabytes of memory. This reliance severely restricts the ability to simulate large quantum systems, hindering scientific discovery.

While recent developments in state-vector simulators have made significant strides in performance improvement [11, 27, 65, 66], optimizing memory usage remains a largely overlooked area. Tensor network simulation is expected to address this issue [38, 42] by representing the quantum circuit using tensor structures and employing tensor contraction to compute the final state vector amplitudes. However, tensor network simulators face significant limitations when simulating highly entangled quantum circuits [41, 57]. For entanglement-heavier circuits, both the computational and memory overhead of tensor network simulators grow substantially. This restricts their applicability primarily to circuits that are shallow and exhibit low entanglement between qubits. For instance, using tensor network simulators to execute the Quantum Approximate Optimization Algorithm (QAOA) [13] and the Variational Quantum Eigensolver (VQE) [43], the most representative quantum algorithms in the NISQ era, faces significant limitations. In QAOA, tensor networks can only efficiently manage a limited number of layers

[37], while an arbitrary number of operational layers is essential to increase effectiveness [12, 20]. For VQE, the enormous number of gates and the level of qubit entanglement [54] create impractical scenarios for tensor networks to solve.

That said, state vector-based quantum circuit simulation offers generality and universal benefits for simulating complex quantum algorithms. To this end, relaxing the memory constraint for state vector simulation is the top priority task. In the classical HPC domain, data compression has proven effective in multiple scientific areas for memory reduction. Broadly speaking, compression techniques can be classified into lossy and lossless, based on the trade-offs between the error and compression ratio they introduce to the data. Compared to lossless compression, lossy compression tends to provide better compression rates [63, 64], making it more suitable for high-memory burden scenarios like quantum simulation. Recent studies [6, 23, 56, 63] have developed error-bounded lossy compressors on GPUs, achieving a balance between compression ratios, high-quality data reconstruction, and performance. Incorporating these advanced compression algorithms into quantum simulation holds considerable promise for significantly reducing memory demands, thereby addressing the fundamental challenge in the field.

However, the direct application of a compression technique on state-vector simulations is inefficient and may result in low simulation fidelity. A prior study [58] introduces a workflow that addresses this integration. The workflow starts with compressing the entire state vector. For each gate in the circuit, it breaks the compressed elements into blocks, decompresses each block, updates the state elements in the block, and then re-compresses it until all blocks are processed. This design introduces several potential complications, particularly concerning the performance of the simulation and the fidelity of the quantum state. These issues encompass five primary domains:

Challenge 0: Frequent Compression. Since the entire state vector needs to be updated when simulating each quantum gate, a large quantum system would require frequent compression and decompression operations on the critical path of the state vector simulation, introducing significant performance overhead. Moreover, lossy compression inherently introduces errors into the reconstructed data. When simulating deep quantum circuits, these errors accumulate and degrade the fidelity of the final results.

Challenge @: Memory Movement Overhead. To maximize the number of qubits supported by simulation and improve the simulation performance, the involvement of large memory space such as CPU memory and high-parallelism computing resources like GPUs is necessary. However, the data movement between the CPU and GPU to take advantage of computation acceleration incurs significant overhead.

Challenge O: Lack of Dedicated Error Control Scheme. Effective error control in lossy compression is essential, particularly for the point-wise relative error control scheme for state-vector simulation [58]. The GPU-based compression processes would outperform their CPU-based alternatives and eliminate potential additional memory transfers between the CPU and GPU. However, current GPU-based lossy compressors do not incorporate such a scheme.

Challenge 4: Unpredictable Memory Consumption of Compressed State Vectors. When handling large input data, lossy compressors often divide the data into smaller chunks for independent compression. However, the memory footprints of the compressed state vector chunks depend on the properties of the state vector, complicating the accurate assessment of whether the available memory will suffice for the simulation.

In response to these challenges, we introduce a novel state vector quantum simulation framework, BMQSIM, by efficiently integrating lossy compression techniques. This framework can break the memory limit to support the robust simulation of more qubits on GPU machines while maintaining high fidelity in simulation results by significantly reducing the frequency of compression with a novel circuit partition scheme. BMQSIM is adaptable, allowing for easy integration into various simulators, enhancing its utility across different simulation backends.

Our paper makes the following contributions:

- We introduce a novel circuit partitioning strategy, effectively addressing low-fidelity and low-performance concerns of the compression-integrated simulation. This method divides the simulation process into discrete subtasks, each involving a partition of the circuit and corresponding elements of the state vector. This approach significantly reduces the frequency of compression and decompression operations, thereby maintaining exceptionally high simulation fidelity and significantly improving simulation time.
- We propose an innovative workflow pipeline that concurrently executes (de)compression operations and data movement. This approach minimizes the perceived overhead in the simulation process by effectively *hiding* these operations within the data transfer time frames.
- We develop the first GPU-based point-wise error control mechanism in a lossy compressor. It offers adaptability to other compressors requiring absolute error control, marking a significant advancement in GPU-accelerated data compression.
- We propose a two-level memory management system to address the challenge of unpredictable compressed state vector block sizes. It dynamically manages memory (de)allocation and uses the GPUDirect Storage technique

to create an effective secondary memory buffer in an SSD, ensuring efficient memory utilization and enhanced operational stability.

• Evaluations on various circuits demonstrate that BMQSIM significantly enhances the capabilities of SOTA state-vector simulators by enabling the simulation of up to 14 additional qubits (on average 10 additional qubits) under the same memory constraints, while maintaining comparable simulation times to SOTA simulators.

This paper is organized as follows: §2 provides background information. §3 analyzes the problem and discusses the issues of basic solutions. §4 details our design. Evaluation results are presented in §5. Finally, §8 summarizes our findings and discusses future research directions.

2 Background

In this section, we introduce state-vector simulation, floatingpoint data compression, and CUDA architecture.

2.1 Principles of State-Vector Simulation

In quantum computing, a qubit, like a bit in classical computing, is the fundamental unit for computing. Unlike bits in traditional computing, a qubit can have many more states besides 0 and 1. A qubit $|\psi\rangle$ is a two-level state that can be expressed as:

$$|\psi\rangle = a_0|0\rangle + a_1|1\rangle$$

Here, a_0 and a_1 represent two complex amplitudes, where $|a_0|^2 + |a_1|^2 = 1$. The quantum state with *n* qubits can be described as a state vector containing 2^n complex amplitudes:

$$|\psi\rangle = a_{0\cdots00}|0\cdots00\rangle + a_{0\cdots01}|0\cdots01\rangle + \cdots + a_{1\cdots11}|1\cdots11\rangle$$

This state also adheres to the condition $\sum_i |a_i|^2 = 1$. The subscripts of *a* are the indices in binary format. In the computation of simulation, the state vector is often denoted as a column vector:

$$\begin{bmatrix} a_{0\cdots 00} \\ a_{0\cdots 01} \\ \vdots \\ a_{1\cdots 11} \end{bmatrix}$$

A quantum gate represents a unitary operation applied to qubit(s), and a series of quantum gates operating on a set of qubits forms a quantum circuit. Applying a gate to a qubit is equivalent to conducting a matrix multiplication of the gate unitary matrix and the elements in the state vector. These matrices modify the elements of the state vector corresponding to the target qubit(s). The most common types of gates are single-qubit gates and double-qubit gates. For a single-qubit gate (a 2×2 matrix) applied to qubit *k*, the operation is to multiply the matrix with two elements whose indices differ only in the *k* bit:

$$\begin{bmatrix} a'_{e_{2^{n}-1}\cdots 0_{k}^{1}\cdots e_{0}^{1}} \\ a'_{e_{2^{n}-1}\cdots 1_{k}^{2}\cdots e_{0}^{2}} \end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \end{bmatrix} \begin{bmatrix} a_{e_{2^{n}-1}^{1}\cdots 0_{k}^{1}\cdots e_{0}^{1}} \\ a_{e_{2^{n}-1}^{2}\cdots 1_{k}^{2}\cdots e_{0}^{2}} \end{bmatrix}$$
$$\forall \begin{bmatrix} a_{e^{1}} \\ a_{e^{2}} \end{bmatrix}, e_{i}^{1} = e_{i}^{2} \quad \text{for } 0 \le i < 2^{n} \text{ and } i \ne k$$

where $[a_*]$ are the state vector amplitudes and $[u_*]$ is the unitary matrix of the applied gate. Similarly, for a doublequbit gate (a 4 × 4 matrix) applied to qubits q and k, the matrix operation is:

$$\begin{bmatrix} a'_{e_{2^{n}-1}\cdots 0_{q}^{1}\cdots 0_{k}^{1}\cdots e_{0}^{1}}\\ a'_{e_{2^{n}-1}\cdots 1_{q}^{2}\cdots 1_{k}^{2}\cdots e_{0}^{2}}\\ a'_{e_{2^{n}-1}\cdots 1_{q}^{1}\cdots 1_{k}^{1}\cdots e_{0}^{1}}\\ a'_{e_{2^{n}-1}\cdots 1_{q}^{1}\cdots 1_{k}^{1}\cdots e_{0}^{1}}\\ a'_{e_{2^{n}-1}\cdots 1_{q}^{1}\cdots 1_{k}^{1}\cdots e_{0}^{1}}\end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} & u_{13} & u_{14}\\ u_{21} & u_{22} & u_{23} & u_{24}\\ u_{31} & u_{32} & u_{33} & u_{34}\\ u_{41} & u_{42} & u_{43} & u_{44} \end{bmatrix} \begin{bmatrix} a_{e_{2^{n}-1}\cdots 0_{q}^{1}\cdots 0_{k}^{1}\cdots e_{0}^{1}}\\ a_{e_{2^{n}-1}\cdots 1_{q}^{1}\cdots 0_{k}^{2}\cdots e_{0}^{2}}\\ a_{e_{2^{n}-1}^{3}\cdots 1_{q}^{3}\cdots 0_{k}^{3}\cdots e_{0}^{3}\\ a_{e_{2^{n}-1}^{4}\cdots 1_{q}^{4}\cdots e_{0}^{4}}\end{bmatrix} \\ \forall \begin{bmatrix} a_{*} \end{bmatrix}, e_{i}^{1} = e_{i}^{2} = e_{i}^{3} = e_{i}^{4} \quad \text{for } 0 \leq i < 2^{n} \text{ and } i \neq k, i \neq q \end{cases}$$

An important requirement for both single-qubit gates and double-qubit gates is that simulating a gate operation requires iterating through the entire state vector.

2.2 Floating-Point Lossy Compression

In the field of data compression, there are two main types: lossless and lossy compression. Lossless compression retains the original data perfectly, while lossy compression, in exchange for a higher compression ratio, incurs some loss of accuracy. The latter is suitable for scenarios where a certain level of data degradation is acceptable.

Recently, there have been significant advancements in lossy compression algorithms, particularly for floating-point scientific data[5, 10]. Prominent examples are SZ [9, 31, 33, 53, 60], ZFP [34], MGARD [17, 32], and TTHRESH [3]. These algorithms are distinct from traditional lossy compressors for images/videos, as they feature precise error-controlling schemes. These schemes allow for control over the level of accuracy in reconstructed data and further data analysis.

With the rise of GPU-based systems, compatible versions of these compressors, such as cuSZ [7, 56], cuSZ-i[35], cuZFP [8], and MGARD-GPU [6], have been developed using CUDA [46]. Furthermore, new GPU-oriented lossy compressors like FZ-GPU [63], bitcomp [39], and cuSZp [22, 23] have emerged. These GPU versions typically offer higher compression throughputs than their CPU counterparts, enabling their application to a wide range of scenarios, such as deep learning training acceleration [14] and communication acceleration [21]. Other platforms have also been explored for similar reasons, such as Cerebras [49, 50].

However, a gap remains in current GPU compressors: most only support absolute error control or fixed-rate modes. The former keeps the maximum error within a user-defined limit, while the latter targets a specific compression ratio. A critical missing feature is a point-wise relative error control scheme, vital for state-vector simulation to ensure high fidelity [58].

2.3 CUDA Memory Architecture

The increasing adoption of GPUs as the main accelerators of high-performance computing tasks is primarily due to their superior parallel computation capabilities. Within the CUDA architecture [46], a widely used programming model for GPUs, processing units are organized into threads. These threads are grouped into blocks and then organized as a grid structure. GPUs typically feature on-chip memory, or device memory, which is usually much less abundant compared to CPU memory or main memory.

Most applications initialize memory allocation on the CPU and then copy the data to the GPU for computation through PCIe. Therefore, asynchronous memory copy operations are crucial for reducing data transfer latency between the CPU and GPU. Such operations enable GPU kernels (GPU processes) to run concurrently with memory copy tasks, optimizing data transfer efficiency. Recently, data copying can occur directly between SSDs and GPU memory. For the movement of data between SSDs and GPUs, the GPUDirect Storage (GDS) technique is vital. This technique allows GPUs to directly access data stored on SSDs, bypassing the CPU thus enhancing performance.

3 Feasibility Analysis

In this section, we provide a detailed analysis of the solution developed in SC19-Sim [58] to integrate data compression with state vector simulation and identify its shortcomings. *For simplicity, from now on, we use single-qubit gates and binary representation of indices in all the following examples.*

The prior work [58] proposed a basic solution of applying compression techniques in state vector simulation. This solution consists of two key designs: state vector partitioning and state vector updating.

State Vector Partitioning. To maximize flexibility and enable parallel execution of compression and simulation, SC19-Sim divides the state vector into blocks, which we term <u>SV blocks</u>. A demonstration of the state vector partition is illustrated in Figure 2. Assume that the state vector is divided into 2^c SV blocks, and each SV block contains 2^b state vector elements (i.e., amplitudes). Given an <u>n</u>-qubit system, where n = b+c, the higher *c* bits in the qubit index space are referred to as the global index, while the lower *b* bits are referred to as the <u>local index</u>. A clear observation is that within each SV block, the global index remains the same, but the local index varies. Different SV blocks have different global indices.

State Vector Updates. At the beginning of the simulation, the state vector (SV) blocks are compressed and stored in the



Figure 1: An example of how SV blocks are involved based on target global index changes. The same alphabet denotes the same 0 or 1.



Figure 2: An illustration of SV partitioning. We refer to the higher *c* bits as the global index and the lower *b* bits as the local index.

system memory. During the simulation, each gate updates the entire state vector once (as discussed in §2.1). This process involves decompressing every SV block, updating the amplitudes within it, and then recompressing it back into the system memory. Depending on whether the target qubit of the quantum gate is located in the global index or the local index, the updating process may involve either two separate SV blocks or a single SV block, as illustrated in Figure 3. We summarize the updating rules as an observation.

Observation: If the target qubit t_i is in the local index set, the amplitudes needed for matrix-vector multiplication are within the same block. Otherwise, the amplitudes are in different SV blocks, where their exact positions depend on the target qubit.

Issues of the Basic Solution. Based on this observation, the order of processing SV blocks in the simulation process may vary due to the order of the different target qubits of the gates in the circuit. Therefore, without careful design, SC19-Sim applies each gate sequentially to the state vector, requiring decompression and compression before and after updating the state vector amplitudes.

This design exposes several issues: • Since (de)compression is executed on a per-gate basis, fully decompressing all SV blocks for every gate operation significantly lowers performance. Moreover, as the circuit length (number of gates) increases, the number of lossy compression operations escalates, leading to an accumulation of errors and degradation of state fidelity. • The GPU is not leveraged, as the entire state vector is processed only by the CPU. Leveraging the parallel computing capability of



Figure 3: A demonstration of how the target qubit location influences amplitude updates. The same alphabet denotes the same value.

GPUs can significantly improve performance. However, the intensive data transfer between CPU and GPU will heavily impact simulation efficiency. ⁽¹⁾ The compression-introduced error is not controlled. Random errors introduced by compression will result in unguaranteed fidelity. Therefore, we need a specialized error control scheme to bound the fidelity. ⁽¹⁾ The compression ratio is unpredictable during the simulation process. The simulation may halt midway due to insufficient memory space, necessitating a backup memory management system to prevent such interruptions.

4 Design of BMQSIM

Overview of BMQSIM's Design. BMQSIM is designed to simulate full-state quantum circuits with a smaller memory footprint to support more qubit systems. Figure 4 summarizes the key techniques implemented in BMQSIM and the respective sections where they are discussed. Specifically, we introduce a specialized circuit partition approach (§4.1) to minimize the (de)compression frequency hence significantly improve the performance and increase the fidelity, addressing **①**. We include a workflow design (§4.2) to overlap compression/decompression, data movement between CPUs and GPUs, and computation, addressing **②**. An error-controlled GPU compressor (§4.3) is proposed to mitigate **③**. Finally, we present a two-level memory management system (§4.4) to solve issue **④**.

4.1 Optimal-Compression Circuit Partition

As analyzed in Section 3, the basic solution will lead to frequent (de)compression because gates in the circuit require different access patterns on SV blocks due to different target qubits. This issue significantly impacts the simulation performance.

Insight from the Analysis. To solve this issue, we carefully analyze Observation in §3 and obtain two important findings. (1) For multiple gates targeting the local index set,



Figure 4: An overview of our proposed BMQSIM.



Figure 5: An example of the proposed circuit partition process.

we can apply them all after decompressing the corresponding SV block because every amplitude in this block can find its corresponding pair within the same block. (2) For multiple gates targeting the global index set, since different gates may require pairs of different SV blocks, we can involve a few more SV blocks to make the multi-gate application possible and balance the far-reach of pairs. An example of this is shown in Figure 1: when two gates targeting different global indices are applied, we can include more SV blocks to ensure that the pairs of amplitudes needing updating can still be found within these SV blocks. The number of SV blocks involved is two to the power of the number of targeted global indices. Insight is drawn from above findings:

Insight: If all the gates in the circuit target the local index or a few specific global indices, then the state vector update can be done for all the gates with one decompression.

How to make the circuit consist only of gates targeting certain indices? We find that if we partition the circuit into multiple stages where the number of global indices targeted by the involved gates in a stage is limited, then within such stage, all the gate operations can be performed using the same SV block access pattern. Therefore, we propose a circuit partition algorithm to partition the circuit into stages given a pre-defined limit of number of global indices. Details of this approach can be found in Algorithm 1. Specifically, we define global indices that appear within a stage as <u>inner indices</u> and other global indices that do not as <u>outer indices</u>. After the user specifies the SV block size and the inner size, the algorithm runs **offline** for a given circuit. For each stage, we add one gate at a time from the input circuit (Line 11) until the number of global indices in the stage reaches a threshold (Lines 7-9). We repeat this process until the circuit is fully traversed (Line 4). Note that the minimum number of inner indices must be two (Line 3). This requirement stems from the structure of quantum circuits, consisting of single- and double-qubit gates. Ensuring at least two inner indices is crucial for effective circuit partitioning when a double-qubit gate's target qubits both fall within the global indices.

An example of this process is depicted in Figure 5. In this example, we partition the circuit into four stages with Algorithm 1. For this 6-qubit (n = 6) circuit, the local index size is 2 (b = 2), and the global index size is 4 (c = 4). In the example stage from the step 4 in Figure 5, indices 3 and 5 are the inner indices of this stage, while 2 and 4 are the outer indices. All the gate operations in this stage only involve the SV blocks with the same outer indices. We call this set of SV blocks an SV group; there are a total of 4 groups in this example. Each group can be updated independently.

Algorithm 1 Proposed circuit partition method.					
Input: circuit, SV block size, inner size					
Output: stages					
1: stages = []					
2: current stage = []					
3: threshold = max(inner size, 2)	▶ 2 for double-qubit gates				
4: while i < number of gates in circuit do					
5: current gate = circuit[i]					
6: query the global indices in [current s	stage + current gate]				
7: if exceed the threshold then	▶ Partition current stage				
8: add current stage to stages					
9: current stage = []	▹ Clear current stage				
10: end if					
11: add current gate to current stage					
12: i++					
13: end while					
14: if current stage not empty then					
15: add current stage to stages					
16: end if					

With this design, each stage requires only one compression and one decompression operation, significantly reducing the frequency of compression. For instance, in the simulation of a 33-qubit QFT circuit, our approach can decrease the number of compression occurrences from 2,673 (i.e., the number of gates) to just 28 (i.e., the number of stages). This substantially increases the final result's fidelity and also improves overall performance.

4.2 Transfer-Concealed Workflow

On one hand, to maximize qubit support, it is beneficial to store compressed state-vector (SV) blocks in the larger CPU memory (e.g., 16GB to 512GB) compared to GPU memory (e.g., 4GB to 80GB). On the other hand, GPU-based simulators outperform CPU-based ones due to high parallelism ideal for matrix multiplication. Therefore, BMQSIM leverages both CPU and GPU: compressing SV blocks in CPU memory and assigning state vector updates to GPUs. This design, however, requires frequent CPU-GPU memory transfers, complicating block-wise state vector updates.

Pipeline design. To resolve this issue, we propose a memory transfer and computation overlapping pipeline. As described in §4.1, the simulation is divided into discrete, independent tasks called SV groups, allowing for more modular and efficient processing. This characteristic is utilized to overlap kernel executions with data transfers (as mentioned in §2.3, GPUs can perform memory copy operations and kernel execution concurrently). A demonstration of this pipeline design is shown in Figure 6: each SV group undergoes a sequence of operations including host-to-device memory copy, decompression, state vector updating, compression, and device-to-host memory copy. These operations are scheduled on the same CUDA stream to maintain the correct execution order. Additionally, operations for different SV groups are scheduled to each CUDA stream repeatedly, facilitating the overlap of overall processes. Moreover, kernel executions can also be overlapped by the GPU scheduler to fully leverage the computing resources in the GPUs. This strategy efficiently overlaps memory operations and kernel execution, enhancing overall performance.

Multi-GPU parallelization. Since the simulation process is divided into independent tasks by our circuit partition, different GPUs can simultaneously process distinct SV groups of SV blocks. This enables native support for concurrency at the inter-GPU level in BMQSIM. As shown in Figure 6, each GPU handles partial SV groups and processes them locally without GPU-to-GPU communication. Note that the throughput of multi-GPU parallelization is bounded by the PCIe bandwidth, as all data transfer between the CPU and GPUs occurs through PCIe. When memory movement is



Figure 6: A demonstration of our multi-stream pipeline design.

intensive, it can cause a starvation problem for GPUs (evaluated in Section 5.8).

Note that in the beginning of the simulation, the state vector is initialized to a standard base state (the first element is 1, all the others are 0) as a common practice [28]. When the initial state differs from this standard as the simulation proceeds, a few quantum gates can be used to establish the desired initial state. After partitioning the state vector, all SV blocks, except the first one, consist only of zeros. Therefore, there is no need to compress the same SV block multiple times. During the initial compression, we only need to compress the block with the first element set to one and another block containing all zeros. Then, we can copy the compressed SV block with all zeros multiple times. This approach reduces the (de)compression overhead by one instance.

4.3 Point-wise Error Control for GPU Compression

We introduce our proposed GPU point-wise compression error control to ensure that the compression-error propagation in simulation can be bounded in the final results.

It has been proven that GPU lossy compression has much better performance and similar compression ratios compared to CPU compression. To this end, we employ GPU lossy compression in BMQSIM to minimize the compression overhead. Previous work has demonstrated a lower bound on the fidelity of the state vector when applying a point-wise relative error bound [58]. Unfortunately, to our knowledge, current SOTA GPU lossy compressors do not support the point-wise relative error bound mode. To address this, we propose the GPU Point-wise Error Compression algorithm. Drawing on previous work by Liang *et al.* [30], we use a logarithmic transformation to convert point-wise relative error bounds to absolute error bounds.

Specifically, let $f(x) = \log_2(x)$ be a bijective transformation of the original data point x. Applying an absolute error bound b_a to f(x) results in the original data being bounded by a point-wise relative error b_r , as shown by the equation:

$$\frac{|f^{-1}(f(x) + b_a) - x|}{|x|} \le b_r \tag{1}$$

The relationship between b_a and b_r can be expressed as:

$$b_a = g(b_r) = \log_2(1+b_r)$$
 (2)

As a result, we can achieve point-wise relative error bounds.

Challenges. Note that the \log_2 transformation in Equation (2) requires positive input values, but satisfying this requirement is a non-trivial task. A common method is to convert negative values to absolute values before applying the \log_2 transformation and use an array to record their indices. However, this approach would significantly lower the overall compression ratio due to the extra space for the index array, potentially even leading to data size inflation.

Our solution. To address this challenge, we propose an algorithm that avoids using an index array to record the negative values. We detail this algorithm in Algorithm 2 (the decompression process is simply the inverse). Specifically, we use a bitmap to store the sign of each number in the original array (Line 1), designating 0 for positive values and zeros (Line 8), and 1 for negative values (Line 5). Then, we convert the negative values to their additive inverse (Line 6) and apply the log transformation (Line 10). Subsequently, we apply lossy compression with absolute error-bounded mode to the data to achieve point-wise error control (Line 15).

Note that based on our observations, bitmaps frequently exhibit long sequences of repeated 0-bits or 1-bits, indicating that the sign of the state vector is often repeated over extensive distances. To address this redundancy, we propose a pre-scan of the bitmap (Line 16). Specifically, the bitmap is partitioned into chunks, within which CUDA's warp-level fast scan functions, __ballot_any and __ballot_all, are employed. These functions, optimized by register direct data exchange, rapidly assess large bitmap chunks to determine if all bits within a chunk are all-0 or all-1. The results are recorded, and redundant all-0 or all-1 chunks are removed. The remaining data is finally compressed using an additional lossless encoding method (Line 17). This approach not only increases the compression ratio but also enhances overall compression performance.

4.4 Two-Level Memory Management

The point-wise error-bounded lossy compression introduced in BMQSIM raises a potential issue: no sufficient memory guarantee for simulation due to variable compression ratios during the simulation. To address this, we propose a twolevel memory management system. Specifically, if the main memory is insufficient, the machine's storage component is employed as a fallback strategy to support the simulation.

Challenges. A couple of reasons make this solution challenging: 1. Data transfer from the storage to the GPU requires an intermediate step of involving CPU memory, needing additional memory space as a temporary buffer for SV blocks from the storage. 2. Moving SV blocks from the storage to

Algorithm 2 GPU point-wise relative error control compression.
Input: SV blocks
Output: compressed SV blocks, compressed bitmap
1: bitmap = []
2: while i < number of SV blocks do > Pipelined in Section 4.2
3: while j < number of elements in SV block do
4: if SV block[i][j] < 0 then
5: add 1-bit to bitmap ► 1 denotes a negative number
6: SV block[i][j] = -SV block[i][j] ▷ Convert to positive
7: else
8: add 0-bit to bitmap ► 0 denotes a non-negative number
9: end if
10: SV block[i][j] = log_2 (SV block[i][j]) > Convert to log scale
11: j++
12: end while
13: i++
14: end while
15: lossy encode (SV block)
16: pre-scan(bitmap)
17: lossless encode (bitmap)

CPU memory and then to GPU memory generates significant latency, degrading the overall simulation performance.

Our solution. To address these challenges, we employ the GDS technology (as introduced in §2.3) to enable direct memory access between GPU and storage, leveraging the Direct Memory Access (DMA) engine. This method bypasses the potential CPU bounce buffer that traditionally is used as an intermediary for transferring memory between the storage and GPU global memory. Utilizing GDS not only conserves CPU memory—heavily employed for storing compressed SV blocks—but also minimizes CPU overhead. This application of GDS in our design thus enhances BMQSIM's capacity to handle larger quantum simulations more robustly.

During the simulation, if BMQSIM detects that there is insufficient memory for an upcoming compressed SV block, it calls the cuFile APIs [16] to directly save this chunk to the storage via GDS. Our evaluation (§ 5) indicates that the performance drop with two-level management is not significant (i.e., 0.7% on average), highlighting our efficient design.

5 Experimental Evaluation

5.1 Experimental Setup

Machines. Due to the administrative privileges required for driver support for the GDS technique, we conduct our evaluation primarily using the following two machines:

<u>Machine 1</u>: A workstation equipped with a 28-core Intel Xeon Gold 6238R CPU at 2.20GHz and two NVIDIA GTX A4000 GPUs (40 SMs, 16 GB each), along with 128 GB DDR4 memory. This workstation runs Ubuntu 20.04.5 and CUDA 12.3.107. It also includes a Samsung 870 EVO MZ-77E4T0E SSD with a capacity of 4 TB and a SATA 6Gb/s interface. The GPUs in this workstation are connected via PCI Express 4.0.

<u>Machine 2</u>: To evaluate multi-GPU performance speedup, we also include a node from an HPC cluster, which includes a 64-core AMD EPYC 7713 CPU at 2.00GHz and four NVIDIA Ampere A100 GPUs (108 SMs, 40GB each). This system has 256 GB DDR4 memory and runs CentOS 7.4 with CUDA 12.2.91. The GPUs are interconnected using NVLink.

Software. We implement BMQSIM based on SV-Sim [27], primarily because SV-Sim (already merged into NWQSim [51]) is an open-source platform with active maintenance. Furthermore, we base our compression on bitcomp from NVCOMP [39], as bitcomp excels among GPU lossy compressors for its exceptional compression throughput and ratio. Bitcomp integrates both lossless mode and lossy mode. We use lossless mode for bitmap and lossy mode for data. We use a point-wise relative error bound of 10⁻³, as this provides a balanced compression ratio and fidelity.

Baselines. We compare BMQSIM with the following baselines: SV-Sim [27], Qiskit-Aer [24], cuQuantum Appliance [40], and HyQuas [65]. Each of these supports GPU-based state-vector simulation. Additionally, we include a comparison with another state-vector simulation work utilizing compression, referred to as SC19-Sim [58]. However, as the implementation of SC19-Sim is not publicly available, we developed a prototype of SC19-Sim with SV-Sim and SZ2 [31, 52]. A detailed comparison is presented in Table 1.

Table 1: Comparison of Different State Vector Simulators

Existing State Vector Simulators	State Vector Location	GPU Updating?	Use Compression?	
Qiskit	CPU+GPU	\checkmark	×	
SV-Sim	GPU	\checkmark	×	
HyQuas	GPU	\checkmark	×	
cuQuantum	GPU	1	×	
SC19-Sim	CPU	×	\checkmark	
BMQSim	CPU	\checkmark	\checkmark	

Benchmark Circuits. We select eight quantum algorithms from NWQBench [28]. This suite includes quantum circuits with qubit numbers ranging from 23 to 33 and gate numbers from 24 to 3010. The selected circuits are cat_state, cc, ising, qft, bv, qsvm, ghz_state, and qaoa.

5.2 Evaluation of Supported Qubit Number

We begin by assessing the maximum supported number of qubits across different simulators on Machine 1, as shown in Table 2. Our evaluations indicate that BMQSIM can support up to 42 qubits, significantly exceeding other counterparts, which support an average of 30 qubits. This capacity goes beyond some entire HPC clusters under normal simulation conditions. Note that with the help of an SSD (assuming SSDs as an extral external storage space), BMQSIM can reach up to

Table 2: Maximum Qubit Numbers for Different Simulators on Ma-chine 1.

Algorithm	Qiskit	cuQuantum	SV-Sim	HyQuas	BMQSIM
cat_state	33	31	26	29	42
сс	30	N/A	26	29	37
ising	33	31	26	29	35
qft	33	31	26	29	36
bv	33	31	26	29	42
qsvm	33	31	26	29	35
ghz	33	31	26	29	42
qaoa	29	31	26	29	35

47 qubits, which is close to the capacity of the Frontier HPC cluster at 48 qubits [2], and 14 more than other simulators. Note that the supported number of qubits varies due to the unpredictable compression ratio.

5.3 Comparison with SC19-Sim

We then compare BMQSIM with another compression-based state-vector simulation, SC19-Sim [58], to demonstrate the high-performance and high-fidelity advantages of our work.

Since SC19-Sim is not open-source, we implemented a prototype based on SV-Sim [27] with the fastest compression technique, solution B, in their paper [58]. For a fair comparison, we implemented both a CPU version as described in the SC19-Sim [58] paper and a GPU version using the same compression technique but utilizing GPUs to update the state vectors. We ran this evaluation on Machine 1.

Simulation Time. We begin by comparing the simulation time. The results are shown in Figure 7. Our findings indicate that BMQSIM outperforms both versions of SC19-Sim



Figure 7: Simulation time of SC19-Sim (CPU/GPU) and BMQSIM.



Figure 8: Fidelity of SC19-Sim and BMQSIM (higher values better fidelity).

under all configurations. The average speedup of BMQSIM compared to SC19-Sim (CPU) and SC19-Sim (GPU) is 1385× and 539×, respectively. This significant performance boost is attributed to the low compression frequency, finely pipelined workflow, and high-performance GPU compression. Note that in some cases, the SC19-Sim CPU version outperforms the SC19-Sim GPU version. This anomaly is due to the basic solution implemented in SC19-Sim that does not overlaps the data transfer and kernel execution. This results a huge overhead in the memory movement between the CPU and GPU. In contrast, our work leverages a pipeline design to minimize the overhead of data transfer and gain significant performance improvement (evaluated in Section 5.6).

Fidelity. Next, we evaluate the fidelity of simulation results. Fidelity is the most important metric for determining the authenticity of final quatum state. It indicates the similarity between the ideal output state and the simulated state, with values ranging from 0 to 1, where higher is better. The fidelity of our simulations is calculated using the equation: *Fidelity* = $|\langle \psi_{ideal} | \psi_{sim} \rangle|$, where ψ_{ideal} is the ideal output state from SV-Sim and ψ_{sim} is the state produced by the tested lossy-compression enabled simulation. Our results show that BMQSIM achieves a fidelity greater than 0.99 across all configurations, which is higher than SC19-Sim, particularly for deep circuits. For instance, BMQSIM achieves 1.35× higher fidelity on average compared to SC19-Sim in the qft circuit.

5.4 Evaluation of Memory Consumption

We present a memory consumption comparison between BMQSIM and the standard for state vector simulation, which is 2^{n+4} bytes, where *n* denotes the number of qubits, as shown

in Figure 9. The (de)compression is performed once for each circuit stage. We consider the maximum memory consumption across all stages in the circuit as the final memory consumption of the simulation. Extremely low memory usage is observed for cat_state, bv, and ghz_state, with average memory reductions of 678.61 times for cat_state, 424.77 times for bv, and 678.52 times for ghz_state. Other circuits also maintain significant memory reductions, averaging 15.50 times for cc and 10.54 times for qft.

Note that in most cases, system memory is sufficient for simulation. Thus, to evaluate the two-level memory management design that uses SSD storage as a backup plan for simulation, we limit the memory space of Machine 1 to 8 GB and run the same evaluation. We find that the SSD is leveraged only when the qubit number is larger than 32 qubits for some circuits. For example, the ising circuit stores 39% and 70% of its SV blocks in the SSD with qubit numbers 32 and 33, respectively.

5.5 Evaluation of Simulation Time

Next, we evaluate the simulation time of BMQSIM compared with other baselines on Machine 1, as shown in Figure 10.

Compared to SV-Sim, BMQSIM offers significant performance improvements. When NVLink is not available, SV-Sim experiences substantial overhead from GPU-to-GPU communication, resulting in the longest simulation times across all settings. In contrast, BMQSIM partitions the circuit into stages, dividing the simulation into independent local jobs on GPUs, which eliminates the GPU-to-GPU communication, resulting in an average performance speedup of 75×.

In most cases, BMQSIM achieves similar simulation times to the Qiskit-Aer GPU simulator. For instance, the simulation time ratio of BMQSIM to Qiskit-Aer is 0.99 and 1.05 for qsvm and qft on average, respectively. This demonstrates that BMQSIM has optimized the simulation process to perform on par with the SOTA GPU simulator from industry. It is important to note that Qiskit-Aer utilizes both CPU and GPU memory for storing the state vector and prioritizes GPU memory based on our evaluation. Consequently, there is a significant drop in performance when the qubit number increases from 30 to 31, as the GPU memory becomes insufficient, causing a fallback to combined memory.

Despite the improvements, both cuQuantum and HyQuas still outperform BMQSIM in most cases. This performance disparity is primarily due to the SV-Sim backend on which BMQSIM is based. HyQuas, with its series of performance optimizations, achieves the best performance among all simulators, being 12× faster than BMQSIM on average. However, this performance comes at the cost of higher memory consumption, which limits HyQuas's supported qubit number compared to other GPU simulators. CuQuantum, tested using





Figure 9: Memory consumption of BMQSIM compared to the memory required for normal state vector simulation.



Figure 10: Simulation time on various quantum circuits and qubit numbers (missing bars indicate memory allocation errors).

the backend integrated in qsim [55], achieves approximately 9× speedup compared to BMQSIM. However, cuQuantum is not an open-source tool and only supports the float32 data type. This inherent characteristic renders it faster than all the other evaluated simulators using float64 data points.

Compared to these well-optimized works, the advantage of BMQSIM lies in its ability to support a considerable larger number of qubits. Given the popularity and acceptance in the community, BMQSIM offers comparable simulation time with industry-level simulators like Qiskit with significantly more supported qubits.

5.6 Compression Overhead Analysis

We further evaluate the compression overhead of our design by comparing it with the version of BMQSIM without compression, as shown in Figure 11. For this evaluation, we use a single A4000 GPU in Machine 1 to reduce the impact of other overhead on the evaluation results. The results illustrate that, thanks to our circuit partitioning and pipeline design optimizations, the compression overhead is minimal compared to the version without compression. Notably, in some cases, BMQSIM even outperforms the no-compression version. This is because, although compression adds overhead to the simulation process, it also reduces memory copy time due to the smaller size of the compressed SV blocks. When the compression ratio is high, as in the cases of the *cat_state*, *bv*, and *ghz* algorithms, the data copy overhead becomes negligible, enabling BMQSIM to outperform the version without compression. Overall, the compression technique contributes positively to the simulation time and leads to a 9% speedup on average. In comparison, compression accounts for approximately 61% on average of the SC19-Sim simulation time, demonstrating that our work significantly lifts compression overhead. Note that we also evaluate the overhead introduced by the logarithmic transformation in our compression design, and it is negligible in both compression and decompression time (less than 5%).

5.7 Pipeline Design Analysis

We also evaluate the impact of different CUDA stream numbers and present the results on Machine 1 in Figure 12. We fix other parameters, such as the SV block size and inner size, to isolate the impact of the stream number. When the CUDA stream number is set to 1, it represents the version of BMQSIM without pipeline optimization. Our findings indicate that, in most cases, the highest speedup is achieved when the stream number is set to 2. Although the speedup is not as significant with a stream number of 4, some improvement is still observed. However, when the stream number reaches 8, the pipeline version becomes slower than the sequential ICS '25, June 08-11, 2025, Salt Lake City, UT, USA





version. This is due to the stream context switch overhead outweighing the benefits brought by pipeline speedup.



Figure 12: Impact of CUDA stream number in our pipeline design.

5.8 Other Evaluations

Finally, we evaluate other settings, including the GPU number, inner size, SV block size, and circuit partition overhead.

Multi-GPU Speedup. To evaluate the scalability of our work, we test it on up to 4 A100 GPUs from Machine 2 with different circuits of 28 qubits, as shown in Figure 13. In the qft, our work achieves a speedup of $1.7 \times$ and $2.3 \times$ for 2 GPUs and 4 GPUs, respectively, thanks to the independent SV groups design in BMQSIM. While the speedup is not significant when the number of GPUs rise from 2 to 4 in some cases due to the CPU and GPU memory transfer rate bounded by the PCIe (as mentioned in Section 4.2) and the high overhead of GPU operation launches.



Figure 13: Scalability of BMQSIM on different algorithms.

Evaluation of Circuit Partition Overhead. To demonstrate the overhead of the extra circuit partition strategy, we evaluate the percentage of the circuit partition time compared to the end-to-end latency of the simulation process, as

shown in Figure 14. The results indicate that the partition time is negligible compared to the overall simulation time.



Figure 14: Circuit partition time as a percentage of overall simulation time.

Parameter Tuning. To evaluate the influences of the inner size and SV block size, we assess the simulation time and compression ratio (the ratio of standard memory to practical memory) with different settings for the 30-qubit qaoa algorithm, as shown in Figure 15. Our findings indicate that the compression ratio does not vary significantly with different inner sizes and SV block sizes. However, the simulation time is shorter with higher inner sizes and SV block sizes. This is because a larger inner size and SV block size result in fewer stages and, consequently, fewer kernel launches.



Figure 15: The impact of two system parameters (i.e., inner size and SV block size) on compression ratio (left) and simulation time (right).

Fidelity Evaluation. To demonstrate the high-fidelity performance of BMQSIM, we include results for larger numbers of qubits, as shown in Table 3. The results show that

Table 3: Algorithm fidelity across qubit counts 24 to 30.

Algorithm	24	25	26	27	28	29	30
cat_state	0.9995	0.9995	0.9995	0.9995	0.9995	0.9995	0.9995
сс	0.9996	0.9998	0.9987	0.9988	0.9997	0.9998	0.9993
ising	0.9992	0.9997	0.9987	0.9993	0.9992	0.9997	0.9998
qft	0.9998	0.9993	0.9988	0.9983	0.9998	0.9993	0.9988
bv	0.9985	0.9995	0.9995	0.9995	0.9995	0.9995	0.9995
qsvm	0.9992	0.9997	0.9987	0.9993	0.9992	0.9997	0.9998
ghz	0.9995	0.9995	0.9995	0.9995	0.9995	0.9995	0.9995
qaoa	0.9986	0.9999	1.0000	0.9994	0.9990	1.0000	0.9998

BMQSIM consistently achieves fidelity above 0.99 across all algorithms and qubit counts, highlighting its robustness and accuracy. It also indicates that BMQSIM is suitable for a wide range of algorithms, with minimal impact from varying entanglement patterns.

6 Discussion

In the previous sections, we introduce BMQSIM and show that it can achieve up to 14 more qubits for state-vector simulation. A follow-up question is can one apply it to densitymatrix and tensor-network simulations? Unlike the statevector approach, the density-matrix and tensor-network approaches represent quantum states using matrices and tensors, respectively. Lossy compression can be applied to arrays of any dimension, including 1D vectors, 2D matrices, and 3D/4D tensors, where typically a higher dimension allows for a higher compression ratio. However, the strategy for updating elements in the density matrix differs from that in the state vector and includes irregular access, which necessitates a different partition strategy. In contrast, computation in the tensor network requires the simultaneous reuse of the same part of a tensor, creating a dependency between compression and tensor updates. Addressing these challenges to enable BMQSIM in these issues remains an area for future work.

Regarding *integration with other state-vector simulators*, BMQSIM is designed to be independent of both the simulator's front and back end. We propose techniques like gate remapping to implement our work as a separate plugin. However, there are still some challenges when integrating with other simulators. For example, in combination with HyQuas, the unique circuit partition pattern of this work must be reconciled with our partition strategy. When integrating with cuQuantum, the closed-source nature of its backend impacts the implementation of our pipeline design, as we cannot specify the execution CUDA stream.

7 Related Work

The field of quantum simulation, particularly based on state vectors, has been a significant area of research in recent years.

Various quantum computing machines provide their own simulators, with notable examples including Qiskit [24], Cirq [18] and cuQuantum SDK [40]. These machines, supported by specialized development teams, prioritize stability and versatility in their simulation tools. In recent years, there has been an extensive body of work in state-vector simulation, such as Atlas[59], QX [26], qHiPSTER [48], IQS [19], HiSVSIM [11], Hyquas [65], UniQ [66], SV-Sim [27]. These simulators primarily focus on improving the simulation performance by enhancing memory locality and communication efficiency.

Atlas [59] abstracts the quantum simulation communication optimization problem as a linear programming problem. By solving this problem offline, Atlas achieves improved simulation performance. However, the offline analysis is time-consuming and can significantly exceed the simulation time, rendering it impractical for small-scale problems.

QX [26] enhances the efficiency of quantum operations through optimization techniques such as instruction-level parallelism (e.g., SSE, AVX, and FMA instructions) and multithreading. It also performs gate-specific optimization by leveraging the reduction of floating-point operations and swap-based implementation.

qHiPSTER [48], designed by Intel, is a distributed quantum simulation system capable of simulating up to 42 qubits using 1,000 nodes. It introduces a methodology wherein half of the required data is communicated to a corresponding peer node during amplitude updating for each gate on high-order qubits. After computation, the results are sent back to the original nodes. IQS [19] is an upgraded version of qHiPSTER, focusing on reducing global communication in distributed simulations through strategic qubit mapping.

Zhang *et al.* developed HyQuas [65], an advanced quantum circuit simulator that automatically selects the most efficient simulation methods for different sections of a quantum circuit, based on their patterns. HyQuas integrates two highly optimized methods and a GPU-centric communication pipelining approach to enhance performance. Building on this, Zhang *et al.* also introduced UniQ [66], a programming model for high-performance and portable state-vector simulation, offering unified application-level and hardwarelevel abstractions.

Li *et al.* developed SV-Sim, a scalable PGAS-based statevector simulator [27]. This simulator employs direct peer access for intra-node communication and SHMEM for internode communication, enhancing simulation efficiency. SV-Sim is adept at abstracting various quantum gates across a range of heterogeneous backends, such as CPUs, GPUs, and Xeon Phi. Its compatibility with higher-level quantum programming environments, including IBM Qiskit, Microsoft Q#, and Google Cirq, adds to its versatility. An extension of SV-Sim, NWQSim [51], integrates it with a density-matrix simulator [29] for advanced capabilities. In addition, Fang *et al.* proposed HiSVSIM [11] which designs efficient hierarchical circuit partition (i.e. acyclic graph partition) to achieve faster simulation.

An earlier study [58] (SC19-Sim) explored the use of data compression techniques to reduce the memory footprint in state-vector simulation, focusing on CPU-only approaches. Comp-QSim simply applies compression to segments of state vectors and decompresses them for updates. However, it does not fully integrate compression into the overall computation workflow, serving primarily as a proof of concept for the potential memory reduction ratio achievable with the proposed compression techniques.

8 Conclusion and Future Work

In this paper, we introduced BMQSIM to address memory limitations in quantum simulation. We propose four designs: Circuit Partition, Workflow Pipeline, Point-wise Relative Error Control and Two-level Memory to employ lossy compression creatively and effectively tackling challenges such as low simulation fidelity and high compression overhead, BMQSIM has successfully enabled the simulation of up to 14 (on average 10) additional qubits under memory constraints, with fidelity over 0.999 in almost all cases.

BMQSIM has undergone multiple iterations and has continuously evolved toward higher performance and scalability [61, 62]. In future work, we plan to integrate BMQSIM with other state-vector simulators, such as cuQuantum and Qiskit, to improve performance and usability. Furthermore, we aim to extend this work to multi-node scenarios for largescale simulation.

Acknowledgments

This research is partially supported by the U.S. Department of Energy (DOE) through the Office of Advanced Scientific Computing Research's "Orchestration for Distributed & Data-Intensive Scientific Exploration" and the "Decentralized Data Mesh for Autonomous Materials Synthesis" AT SCALE LDRD at Pacific Northwest National Laboratory. PNNL is operated by Battelle for the DOE under Contract DE-AC05-76RL01830. This work was also supported in part by the National Science Foundation under Grant Numbers 2311876, 2326495, 2247060, and 2247080.

References

- Alán Aspuru-Guzik, Anthony D Dutoi, Peter J Love, and Martin Head-Gordon. 2005. Simulated quantum computation of molecular energies. *Science* 309, 5741 (2005), 1704–1707.
- [2] Scott Atchley, Christopher Zimmer, John Lange, David Bernholdt, Veronica Melesse Vergara, Thomas Beck, Michael Brim, Reuben Budiardja, Sunita Chandrasekaran, Markus Eisenbach, et al. 2023. Frontier:

Zhang et al.

Exploring Exascale. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. 1–16.

- [3] Rafael Ballester-Ripoll, Peter Lindstrom, and Renato Pajarola. 2019. TTHRESH: Tensor compression for multidimensional visual data. *IEEE transactions on visualization and computer graphics* 26, 9 (2019), 2891–2903.
- [4] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd. 2017. Quantum machine learning. *Nature* 549, 7671 (2017), 195–202.
- [5] Franck Cappello, Mario Acosta, Emmanuel Agullo, Hartwig Anzt, Jon Calhoun, Sheng Di, Luc Giraud, Thomas Grützmacher, Sian Jin, Kentaro Sano, et al. 2025. Multifacets of lossy compression for scientific data in the Joint-Laboratory of Extreme Scale Computing. *Future Generation Computer Systems* 163 (2025), 107323.
- [6] Jieyang Chen, Lipeng Wan, Xin Liang, Ben Whitney, Qing Liu, David Pugmire, Nicholas Thompson, Jong Youl Choi, Matthew Wolf, Todd Munson, Ian Foster, and Scott Klasky. 2021. Accelerating multigridbased hierarchical scientific data refactoring on gpus. In 2021 IEEE International Parallel and Distributed Processing Symposium. IEEE, 859– 868.
- [7] cuSZ. [n. d.]. https://github.com/szcompressor/cuSZ. Online.
- [8] cuZFP. [n.d.]. https://github.com/LLNL/zfp/tree/develop/src/cuda_zfp. Online.
- [9] Sheng Di and Franck Cappello. 2016. Fast error-bounded lossy HPC data compression with SZ. In 2016 IEEE International Parallel and Distributed Processing Symposium. IEEE, Chicago, IL, USA, 730–739.
- [10] Sheng Di, Jinyang Liu, Kai Zhao, Xin Liang, Robert Underwood, Zhaorui Zhang, Milan Shah, Yafan Huang, Jiajun Huang, Xiaodong Yu, et al. 2024. A survey on error-bounded lossy compression for scientific datasets. arXiv preprint arXiv:2404.02840 (2024).
- [11] Bo Fang, M. Yusuf Özkaya, Ang Li, Ümit V. Çatalyürek, and Sriram Krishnamoorthy. 2022. Efficient Hierarchical State Vector Simulation of Quantum Circuits via Acyclic Graph Partitioning. In *CLUSTER*. 289– 300. doi:10.1109/CLUSTER51413.2022.00041
- [12] Edward Farhi, David Gamarnik, and Sam Gutmann. 2020. The quantum approximate optimization algorithm needs to see the whole graph: A typical case. arXiv preprint arXiv:2004.09002 (2020).
- [13] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. 2014. A quantum approximate optimization algorithm. arXiv preprint arXiv:1411.4028 (2014).
- [14] Hao Feng, Boyuan Zhang, Fanjiang Ye, Min Si, Ching-Hsiang Chu, Jiannan Tian, Chunxing Yin, Summer Deng, Yuchen Hao, Pavan Balaji, et al. 2024. Accelerating Communication in Deep Learning Recommendation Model Training with Dual-Level Adaptive Lossy Compression. In SC24: International Conference for High Performance Computing, Networking, Storage and Analysis. IEEE, 1–16.
- [15] Jay Gambetta. 2020. IBM's roadmap for scaling quantum technology. https://www.ibm.com/quantum/blog/ibm-quantum-roadmap? mhsrc=ibmsearch_a&mhq=condor.
- [16] GDS cuFile API Reference. [n. d.]. https://docs.nvidia.com/gpudirectstorage/api-reference-guide/index.html.
- [17] Qian Gong, Jieyang Chen, Ben Whitney, Xin Liang, Viktor Reshniak, Tania Banerjee, Jaemoon Lee, Anand Rangarajan, Lipeng Wan, Nicolas Vidal, et al. 2023. MGARD: A multigrid framework for highperformance, error-controlled data compression and refactoring. *SoftwareX* 24 (2023), 101590.
- [18] Google. [n. d.]. Cirq. https://quantumai.google/cirq. Online.
- [19] Gian Giacomo Guerreschi, Justin Hogaboam, Fabio Baruffa, and Nicolas PD Sawaya. 2020. Intel Quantum Simulator: A cloud-ready highperformance simulator of quantum circuits. *Quantum Science and Technology* 5, 3 (2020), 034007.

- [20] Gian Giacomo Guerreschi and Anne Y Matsuura. 2019. QAOA for Max-Cut requires hundreds of qubits for quantum speed-up. *Scientific reports* 9, 1 (2019), 1–7.
- [21] Jiajun Huang, Sheng Di, Xiaodong Yu, Yujia Zhai, Jinyang Liu, Yafan Huang, Ken Raffenetti, Hui Zhou, Kai Zhao, Xiaoyi Lu, et al. 2024. gzccl: Compression-accelerated collective communication framework for gpu clusters. In *Proceedings of the 38th ACM International Conference on Supercomputing*. 437–448.
- [22] Yafan Huang, Sheng Di, Guanpeng Li, and Franck Cappello. 2024. cuSZp2: A GPU Lossy Compressor with Extreme Throughput and Optimized Compression Ratio. In SC24: International Conference for High Performance Computing, Networking, Storage and Analysis. IEEE, 1–18.
- [23] Yafan Huang, Sheng Di, Xiaodong Yu, Guanpeng Li, and Franck Cappello. 2023. cuSZp: An Ultra-fast GPU Error-bounded Lossy Compression Framework with Optimized End-to-End Performance. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. 1–13.
- [24] IBM. [n.d.]. Qiskit. https://qiskit.org/. Online.
- [25] Tyson Jones, Anna Brown, Ian Bush, and Simon C Benjamin. 2019. QuEST and high performance simulation of quantum computers. *Scientific reports* 9, 1 (2019), 1–11.
- [26] Nader Khammassi, Imran Ashraf, Xiang Fu, Carmen G Almudever, and Koen Bertels. 2017. QX: A high-performance quantum computer simulation platform. In *Design, Automation & Test in Europe Conference & Exhibition (DATE), 2017.* IEEE, 464–469.
- [27] Ang Li, Bo Fang, Christopher Granade, Guen Prawiroatmodjo, Bettina Heim, Martin Roetteler, and Sriram Krishnamoorthy. 2021. SV-sim: scalable PGAS-based state vector simulation of quantum circuits. In *SC21*. 1–14.
- [28] Ang Li, Samuel Stein, Sriram Krishnamoorthy, and James Ang. 2021. QASMBench: A Low-level QASM Benchmark Suite for NISQ Evaluation and Simulation. arXiv preprint arXiv:2005.13018 (2021).
- [29] Ang Li, Omer Subasi, Xiu Yang, and Sriram Krishnamoorthy. 2020. Density matrix quantum circuit simulation via the BSP machine on modern GPU clusters. In Sc20: international conference for high performance computing, networking, storage and analysis. IEEE, 1–15.
- [30] Xin Liang, Sheng Di, Dingwen Tao, Zizhong Chen, and Franck Cappello. 2018. An efficient transformation scheme for lossy data compression with point-wise relative error bound. In 2018 IEEE International Conference on Cluster Computing (CLUSTER). IEEE, 179–189.
- [31] Xin Liang, Sheng Di, Dingwen Tao, Sihuan Li, Shaomeng Li, Hanqi Guo, Zizhong Chen, and Franck Cappello. 2018. Error-controlled lossy compression optimized for high compression ratios of scientific datasets. In 2018 IEEE International Conference on Big Data. IEEE, 438– 447.
- [32] Xin Liang, Ben Whitney, Jieyang Chen, Lipeng Wan, Qing Liu, Dingwen Tao, James Kress, David Pugmire, Matthew Wolf, Norbert Podhorszki, et al. 2021. MGARD+: Optimizing multilevel methods for error-bounded scientific data reduction. *IEEE Trans. Comput.* 71, 7 (2021), 1522–1536.
- [33] Xin Liang, Kai Zhao, Sheng Di, Sihuan Li, Robert Underwood, Ali M Gok, Jiannan Tian, Junjing Deng, Jon C Calhoun, Dingwen Tao, et al. 2022. Sz3: A modular framework for composing prediction-based error-bounded lossy compressors. *IEEE Transactions on Big Data* 9, 2 (2022), 485–498.
- [34] Peter Lindstrom. 2014. Fixed-rate compressed floating-point arrays. IEEE Transactions on Visualization and Computer Graphics 20, 12 (2014), 2674–2683.
- [35] Jinyang Liu, Jiannan Tian, Shixun Wu, Sheng Di, Boyuan Zhang, Robert Underwood, Yafan Huang, Jiajun Huang, Kai Zhao, Guanpeng Li, et al. 2024. CUSZ-i: High-Ratio Scientific Lossy Compression on

GPUs with Optimized Multi-Level Interpolation. In *SC24: International Conference for High Performance Computing, Networking, Storage and Analysis.* IEEE, 1–15.

- [36] Seth Lloyd, Masoud Mohseni, and Patrick Rebentrost. 2013. Quantum algorithms for supervised and unsupervised machine learning. arXiv preprint arXiv:1307.0411 (2013).
- [37] Danylo Lykov, Angela Chen, Huaxuan Chen, Kristopher Keipert, Zheng Zhang, Tom Gibbs, and Yuri Alexeev. 2021. Performance evaluation and acceleration of the QTensor quantum circuit simulator on GPUs. In 2021 IEEE/ACM Second International Workshop on Quantum Computing Software (QCS). IEEE, 27–34.
- [38] Igor L Markov and Yaoyun Shi. 2008. Simulating quantum computation by contracting tensor networks. SIAM J. Comput. 38, 3 (2008), 963–981.
- [39] nvCOMP: A library for fast lossless compression/decompression on the GPU. [n. d.]. https://github.com/NVIDIA/nvcomp.
- [40] NVIDIA. [n. d.]. cuQuantum: Accelerate quantum computing research. https://developer.nvidia.com/cuquantum-sdk. Online.
- [41] Stellan Östlund and Stefan Rommer. 1995. Thermodynamic limit of density matrix renormalization. *Physical review letters* 75, 19 (1995), 3537.
- [42] Yuchen Pang, Tianyi Hao, Annika Dugad, Yiqing Zhou, and Edgar Solomonik. 2020. Efficient 2D tensor network simulation of quantum systems. In SC20: International Conference for High Performance Computing, Networking, Storage and Analysis. IEEE, 1–14.
- [43] Alberto Peruzzo, Jarrod McClean, Peter Shadbolt, Man-Hong Yung, Xiao-Qi Zhou, Peter J Love, Alán Aspuru-Guzik, and Jeremy L O'brien. 2014. A variational eigenvalue solver on a photonic quantum processor. *Nature communications* 5, 1 (2014), 4213.
- [44] John Preskill. 2018. Quantum computing in the NISQ era and beyond. *Quantum* 2 (2018), 79.
- [45] Patrick Rebentrost, Brajesh Gupt, and Thomas R Bromley. 2018. Quantum computational finance: Monte Carlo pricing of financial derivatives. *Physical Review A* 98, 2 (2018), 022321.
- [46] Jason Sanders and Edward Kandrot. 2010. CUDA by example: an introduction to general-purpose GPU programming. Addison-Wesley Professional.
- [47] Maria Schuld and Nathan Killoran. 2019. Quantum machine learning in feature hilbert spaces. *Physical review letters* 122, 4 (2019), 040504.
- [48] Mikhail Smelyanskiy, Nicolas PD Sawaya, and Alán Aspuru-Guzik. 2016. qHiPSTER: The quantum high performance software testing environment. arXiv preprint arXiv:1601.07195 (2016).
- [49] Shihui Song, Yafan Huang, Peng Jiang, Xiaodong Yu, Weijian Zheng, Sheng Di, Qinglei Cao, Yunhe Feng, Zhen Xie, and Franck Cappello. 2024. Ceresz: Enabling and scaling error-bounded lossy compression on cerebras cs-2. In Proceedings of the 33rd International Symposium on High-Performance Parallel and Distributed Computing. 309–321.
- [50] Shihui Song, Robert Underwood, Sheng Di, Yafan Huang, Peng Jiang, and Franck Cappello. 2025. A Memory-efficient and Computationbalanced Lossy Compressor on Wafer-Scale Engine. In 2025 ieee international parallel and distributed processing symposium (ipdps). IEEE.
- [51] In-Saeng Suh and Ang Li. [n. d.]. Simulating Quantum Systems with NWQ-Sim on HPC. https://sc23.supercomputing.org/proceedings/ tech_poster/poster_files/rpost195s3-file3.pdf. Online.
- [52] szcompressor. [n. d.]. SZ2. https://github.com/szcompressor/SZ
- [53] Dingwen Tao, Sheng Di, Zizhong Chen, and Franck Cappello. 2017. Significantly improving lossy compression for scientific data sets based on multidimensional prediction and error-controlled quantization. In 2017 IEEE International Parallel and Distributed Processing Symposium. IEEE, Orlando, FL, USA, 1129–1139.
- [54] Andrew G Taube and Rodney J Bartlett. 2006. New perspectives on unitary coupled-cluster theory. *International journal of quantum chemistry* 106, 15 (2006), 3393–3401.

- [55] Quantum AI team and collaborators. 2020. qsim. doi:10.5281/zenodo. 4023103
- [56] Jiannan Tian, Sheng Di, Kai Zhao, Cody Rivera, Megan Hickman Fulp, Robert Underwood, Sian Jin, Xin Liang, Jon Calhoun, Dingwen Tao, et al. 2020. Cusz: An efficient gpu-based error-bounded lossy compression framework for scientific data. In *Proceedings of the ACM International Conference on Parallel Architectures and Compilation Techniques*. 3–15.
- [57] Steven R White. 1992. Density matrix formulation for quantum renormalization groups. *Physical review letters* 69, 19 (1992), 2863.
- [58] Xin-Chuan Wu, Sheng Di, Emma Maitreyee Dasgupta, Franck Cappello, Hal Finkel, Yuri Alexeev, and Frederic T Chong. 2019. Full-state quantum circuit simulation by using data compression. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. 1–24.
- [59] Mingkuan Xu, Shiyi Cao, Xupeng Miao, Umut A Acar, and Zhihao Jia. 2024. Atlas: Hierarchical partitioning for quantum circuit simulation on gpus (extended version). arXiv preprint arXiv:2408.09055 (2024).
- [60] Zhuoxun Yang, Sheng Di, Longtao Zhang, Ruoyu Li, Ximiao Li, Jiajun Huang, Jinyang Liu, Franck Cappello, and Kai Zhao. 2025. PSZ: Enhancing the SZ Scientific Lossy Compressor With Progressive Data Retrieval. arXiv preprint arXiv:2502.04093 (2025).
- [61] Boyuan Zhang, Bo Fang, Qiang Guan, Ang Li, and Dingwen Tao. 2023. Hq-sim: High-performance state vector simulation of quantum circuits

on heterogeneous hpc systems. In Proceedings of the 2023 International Workshop on Quantum Classical Cooperative. 1–4.

- [62] Boyuan Zhang, Bo Fang, Qiang Guan, Ang Li, and Dingwen Tao. 2023. MEMQSim: Highly Memory-Efficient and Modularized Quantum State-Vector Simulation. In Proceedings of the SC'23 Workshops of the International Conference on High Performance Computing, Network, Storage, and Analysis. 1452–1453.
- [63] Boyuan Zhang, Jiannan Tian, Sheng Di, Xiaodong Yu, Yunhe Feng, Xin Liang, Dingwen Tao, and Franck Cappello. 2023. Fz-gpu: A fast and high-ratio lossy compressor for scientific computing applications on gpus. In Proceedings of the 32nd International Symposium on High-Performance Parallel and Distributed Computing. 129–142.
- [64] Boyuan Zhang, Jiannan Tian, Sheng Di, Xiaodong Yu, Martin Swany, Dingwen Tao, and Franck Cappello. 2023. GPULZ: Optimizing LZSS Lossless Compression for Multi-byte Data on Modern GPUs. arXiv preprint arXiv:2304.07342 (2023).
- [65] Chen Zhang, Zeyu Song, Haojie Wang, Kaiyuan Rong, and Jidong Zhai. 2021. HyQuas: hybrid partitioner based quantum circuit simulation system on GPU. In *Proceedings of the ACM International Conference on Supercomputing*. 443–454.
- [66] Chen Zhang, Haojie Wang, Zixuan Ma, Lei Xie, Zeyu Song, and Jidong Zhai. 2022. UniQ: a unified programming model for efficient quantum circuit simulation. In SC22. IEEE Computer Society, 692–707.