



Article

---

# Memory Management Strategies for Software Quantum Simulators

---

Gilberto Díaz, Luiz Steffenel, Carlos Barrios and Jean Couturier

Special Issue

Quantum Computing: A Taxonomy, Systematic Review, and Future Directions

Edited by  
Dr. Yousef Fazea



Article

# Memory Management Strategies for Software Quantum Simulators

Gilberto Díaz <sup>1,\*</sup> , Luiz Steffanel <sup>2</sup> , Carlos Barrios <sup>1,\*</sup>  and Jean Couturier <sup>2</sup> 

<sup>1</sup> Computación Avanzada y a Gran Escala (CAGE), Universidad Industrial de Santander, Bucaramanga 680011, Colombia

<sup>2</sup> CEA, LRC DIGIT, LICIS, Université de Reims Champagne-Ardenne, 51100 Reims, France; luiz-angelo.steffanel@univ-reims.fr (L.S.); jean-francois.couturier@univ-reims.fr (J.C.)

\* Correspondence: gjdiaz@uis.edu.co (G.D.); cbarrios@uis.edu.co (C.B.)

## Abstract

Software quantum simulators are essential tools for designing and testing quantum algorithms on classical computing architectures, especially given the current limitations of physical quantum hardware. This work focuses on studying and evaluating memory management strategies for scalable quantum state simulation. We examine full-state representation, dynamic state pruning, shared-memory parallelization with OpenMP, distributed memory execution using MPI, and error-bounded floating-point compression with ZFP. These techniques are implemented in a prototype simulator and assessed using the quantum Fourier transform as a benchmark, with performance compared against leading open-source simulators such as Intel-QS, QuEST, and qsim. The results show the trade-offs between computational overhead and memory efficiency, and demonstrate that hybrid approaches combining distributed memory and compression can significantly extend the number of qubits that can be simulated. This work contributes practical insights for improving the scalability of software quantum simulators on classical hardware through optimized memory usage.

**Keywords:** quantum computing; high-performance computing; parallel computing



Academic Editors: Vladimir M. Stojanović, Mingxing Luo and Lev Vaidman

Received: 12 June 2025

Revised: 28 August 2025

Accepted: 28 August 2025

Published: 9 September 2025

**Citation:** Díaz, G.; Steffanel, L.; Barrios, C.; Couturier, J. Memory Management Strategies for Software Quantum Simulators. *Quantum Rep.* **2025**, *7*, 41. <https://doi.org/10.3390/quantum7030041>

**Copyright:** © 2025 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

One of the primary reasons to develop quantum computing is that, theoretically, it has been demonstrated that it allows efficient solutions to some complex problems whose best-known solution has an exponential cost for the input size. Quantum superposition, quantum uncertainty, and quantum entanglement are powerful resources that we can use to encode, decode, transmit, and process information in a highly efficient way that is impossible in the classical world.

Recent technological advances have enabled the development of real quantum devices accessed through the cloud. However, these devices are expected to be limited in the short term in terms of the number and quality of their fundamental component, the qubit. Most current quantum devices have a limited number of qubits. In the quantum circuit model, Atom Computing has 1180 qubits, and IBM Osprey has 433 qubits. In the adiabatic model, the D-Wave 2000Q has 2000 qubits. These quantum computers represent prototypes that are not scalable and sufficient to test complex quantum algorithms. Constructing a full-scale quantum computer comprising millions of qubits is a longer-term prospect.

The growing interest in quantum computing and the limitations of real quantum devices have caused many organizations to focus on developing software quantum simulators that run on classical computers. These simulators are trendy tools suitable for testing quantum computing concepts under ideal conditions, avoiding hardware challenges like the limited number and quality of physical qubits and quantum error correction. A list of the very recent initiatives is maintained on several websites [1–4]. This large number of projects reflects the rapid growth of the field and makes it difficult for researchers to decide which tool to use in their research.

Quantum computing simulators, which operate on classical computers, have emerged as valuable and widely used tools in the field of quantum computing. These simulators play a crucial role in the development, testing, and validation of quantum algorithms before they are implemented on actual quantum hardware. One of the primary advantages of quantum simulators is their accessibility. Unlike quantum computers, which are still relatively scarce and often require significant resources and expertise to operate, simulators can be run on conventional computers. This accessibility allows a broader range of researchers and developers to explore quantum algorithms and concepts without the need for physical quantum computing resources.

Quantum simulators offer a controlled environment for designing and refining quantum algorithms. They can simulate ideal quantum systems without the noise and error rates present in current quantum hardware, providing clearer insights into the theoretical performance of an algorithm. This is particularly useful for educational purposes and theoretical research, where understanding the principles of quantum computation and algorithm design is the main focus.

However, the simulation of quantum computing models in classical computers requires exponential time and involves highly complex memory management. The problem is that using conventional techniques to simulate an arbitrary quantum process significantly more prominent than any of the existing quantum prototypes would soon require considerable memory on a classical computer. For instance, to simulate a 60-qubit quantum state, the process would take 18,000 Petabytes or 18 Exabytes of classical computer memory. Therefore, researchers try to reduce such challenges by proposing efficient simulators.

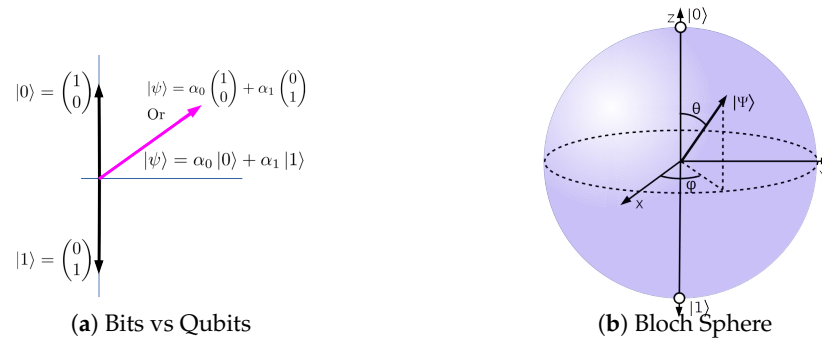
This work explores different aspects involved in the construction of a software quantum simulator capable of performing simulations on classical computers. This includes analyzing and implementing different strategies for memory management, exploring data structures to represent quantum states efficiently, evaluating parallelization techniques such as OpenMP and MPI, and testing compression and pruning methods to optimize performance. Rather than focusing on a single algorithm or optimization, this work aims to provide a comprehensive view of the design decisions, trade-offs, and technical challenges that arise when building scalable quantum simulators on classical hardware.

## 2. Fundamental Concepts of Quantum Computing

To better understand the quantum computing model, it is necessary to know the key principles inherited from quantum mechanics. This section describes the fundamental concepts on which quantum computing is based. Readers already familiar with the field may wish to skip this section.

A logical qubit is a unit vector in a two-dimensional Hilbert space. The Boolean states 0 and 1 are represented by a prescribed pair of normalized and mutually orthogonal quantum states denoted using Dirac's notation  $|0\rangle$  and  $|1\rangle$  [5]. These two states form a *computational basis*, and any other pure state of the qubit can be written as a superposition  $\alpha|0\rangle + \beta|1\rangle$ .

In the classical case, the state of a bit is a special case of a two-dimensional vector: there are only two meaningful orthogonal vectors in the space,  $|0\rangle$  and  $|1\rangle$ , as shown in Figure 1a. In contrast, a qubit is not restricted to these basis states. The general state of a qubit is  $|\psi\rangle = \alpha_0|0\rangle + \alpha_1|1\rangle$ , where  $\alpha_0$  and  $\alpha_1$  are complex numbers satisfying the normalization condition  $|\alpha_0|^2 + |\alpha_1|^2 = 1$ .



**Figure 1.** Visual representation of qubit: a qubit can be written as a superposition  $\alpha_0|0\rangle + \alpha_1|1\rangle$ .

The Bloch sphere is a common representation of a qubit, shown in Figure 1b. Two angles,  $0 < \theta < \pi$  and  $0 \leq \phi \leq 2\pi$ , define the state:

$$|\psi\rangle = e^{i\gamma} \left( \cos\left(\frac{\theta}{2}\right)|0\rangle + e^{i\phi} \sin\left(\frac{\theta}{2}\right)|1\rangle \right).$$

The vector from the origin to the point representing the state makes an angle of  $\theta$  with the z-axis, and its component in the x-y plane makes an angle of  $\phi$  with the x-axis.  $\gamma$  is the global phase, which does not affect the measurable probabilities of the quantum state (it introduces only a uniform Phase Shift to the entire state, which is not physically observable). The state  $|0\rangle$  is the North Pole of the sphere, and the state  $|1\rangle$  is the South Pole.

The general expression of a  $n$ -qubit state is  $|\psi\rangle = \sum_{x=0}^{2^N-1} \alpha_x |X\rangle$  or, in its expanded form,

$$|\psi\rangle = \alpha_0|0\dots 00\rangle + \alpha_1|0\dots 01\rangle + \dots + \alpha_{2^n-1}|1\dots 11\rangle \tag{1}$$

where  $|0\dots 00\rangle = |0\rangle \otimes \dots \otimes |0\rangle \otimes |0\rangle$ ,  $|0\dots 01\rangle = |0\rangle \otimes \dots \otimes |0\rangle \otimes |1\rangle$ , and so on. As we can see, a single complex number can specify a single-qubit state, so  $n$  complex numbers can specify any tensor product of  $n$  individual single-qubit states.

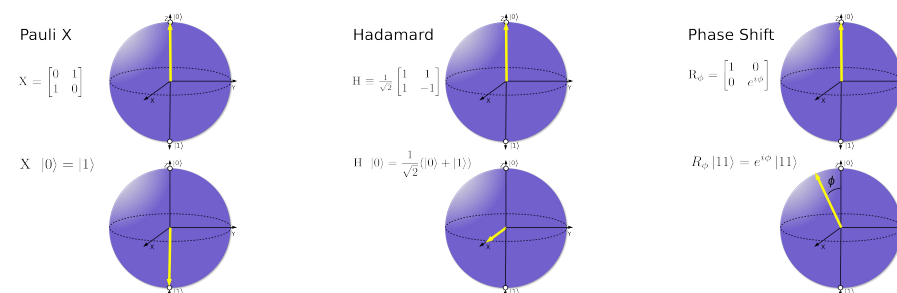
A key feature of quantum mechanics is superposition, where a qubit can exist in a linear combination of  $|0\rangle$  and  $|1\rangle$  until measurement collapses it into one of the basis states, yielding at most one classical bit of information [5]. This restriction on extractable information is formalized by the *Holevo bound*, which establishes an upper limit on the amount of classical information that can be obtained from quantum states [6]. Quantum bits are not restricted to being entirely 0 or 1 in a given instant. In quantum physics, if a quantum system can be found to be in one of a discrete set of states, which we will write as  $|0\rangle$  or  $|1\rangle$ , then, whenever it is not observed, it may also exist in a superposition, or a blend of those states simultaneously [7].

Because a qubit can take on any one of infinitely many states, one can think that a single qubit could store a lot of classical information. However, the properties of quantum measurement severely restrict the amount of information that can be extracted from a qubit. Information about a quantum bit can be obtained only by measurement, and any measurement results in one of only two states, the two basis states associated with the measuring device; thus, a single measurement yields, at most, a single classical bit of information [8].

The quantum entanglement describes a correlation between different parts of a quantum system that cannot be explained classically. It occurs when the subsystems interact so that the resulting state of the whole system cannot be expressed as the direct product of the states of its parts [9,10]. States that cannot be written as the tensor product of  $n$  single-qubit states are called entangled states. Thus, most quantum states are entangled [11]. If we can write the tensor product of those states, they are said to be separable states.

In the quantum circuit model, quantum gates perform unitary transformations on qubits. They are analogous to classical logic gates but operate on quantum states. To transform an  $n$ -qubit state vector, we use  $2^n \times 2^n$  unitary matrices.

Applying a single-qubit gate  $G$  to the  $i$ th qubit corresponds to multiplying the state vector by  $\mathbb{1}_2 \otimes \cdots \otimes \mathbb{1}_2 \otimes G \otimes \mathbb{1}_2 \otimes \cdots \otimes \mathbb{1}_2$ , where  $\mathbb{1}_2$  denotes the  $2 \times 2$  identity matrix. Figure 2 shows some emblematic gates, including the Pauli-X, Hadamard, and Phase Shift gates.



**Figure 2.** Quantum gates: The Pauli X gate acts linearly and it takes the state  $\alpha|0\rangle + \beta|1\rangle$  to the corresponding state in which the role of  $|0\rangle$  and  $|1\rangle$  have been interchanged; it is the quantum equivalent of the NOT gate for classical computers. The Hadamard gate is the first authentic quantum gate because it can generate superposition states. The Phase Shift gate is a single-qubit gate that leaves the basis state  $|0\rangle$  unchanged and maps the state  $|1\rangle$  to  $e^{i\phi}|1\rangle$ .

A quantum algorithm is a step-by-step procedure executed on a quantum computer that exploits quantum features such as superposition and entanglement. While any classical algorithm can be run on a quantum computer, the term “quantum algorithm” usually refers to algorithms that achieve a speedup over classical counterparts, e.g., those based on the quantum Fourier transform, Grover’s search, or quantum simulation.

### 3. Leading Open-Source Simulators for Quantum Computing

Recent years have seen notable advances in the design of software quantum simulators, particularly in scaling to larger qubit counts through improved memory handling and parallelization. Häner et al. [12] explored low-depth tensor network techniques for large-scale circuit simulation, enabling efficient use of memory in circuits with limited entanglement. Villalonga et al. [13] demonstrated a 281PFlop/s simulation of Google’s quantum supremacy circuits, combining distributed memory computing with communication-aware optimizations. Ding et al. [14] proposed a sparse state representation for circuits with low state occupancy, reducing the memory footprint in certain algorithmic classes. Wu et al. [15] introduced GPU-accelerated state-vector simulation strategies that reduce both memory access latency and total storage requirements. These approaches complement the methods investigated in this work, which focuses on full-state simulation with memory optimization through compression, pruning, and hybrid parallelization.

From a recent review, we take some quantum simulators that currently lead the field to characterize their main properties, performance, execution mode, and simulation results to provide comparison and analysis. To facilitate this task, we work with open-source simulators. These simulators are considered state-of-the-art due to several factors [16]:

- **Innovative Features:** Each simulator offers unique capabilities that set them apart, such as optimized algorithms, integration with widely used programming frameworks, or novel approaches to handling quantum state representations. For example, qsim's integration with Cirq and its ability to simulate up to 40 qubits on a high-performance workstation make it a significant tool for developers and researchers.
- **Adoption and Partnerships:** Some of these simulators are backed by major tech companies and have extensive partnerships within the industry, increasing their influence and credibility.
- **Academic and Commercial Use:** These tools are not only used in academic research but are also increasingly adopted by industries for practical applications, which demonstrates their effectiveness and robustness.
- **Recent Updates and Community Support:** The continual updates, community support, and documentation available for these tools contribute to their status as leaders in the field. This ongoing development ensures they remain relevant and useful as quantum computing technology evolves.
- **Open Collaboration:** Open-source projects encourage open collaboration among developers, researchers, and users. Ensuring the source code is available for modification and redistribution fosters a community-driven development approach. This can lead to rapid improvements and innovations, as a diverse group of contributors can work on the software.

The combination of these factors makes these simulators outstanding in the current world of quantum computing, pointing towards their innovativeness and leadership in technological advancement.

### 3.1. Algorithm Selected for Testing

The quantum Fourier transform was selected to carry out the simulations because it offers several advantages: it is a well-studied quantum algorithm with known properties, making it a reliable benchmark for validating the accuracy and efficiency of quantum simulators on classical hardware. QFT's performance scales predictably with the number of qubits, allowing researchers to analyze how the simulator handles increasing complexity. Performing QFT simulations helps estimate the computational resources (memory, processing power) required for larger, more complex quantum algorithms. Finally, QFT is a crucial component in many quantum algorithms, such as Shor's algorithm for factoring large integers. Simulating QFT provides a foundation for testing and understanding these more complex algorithms.

### 3.2. Test Platform

To evaluate the performance of the selected simulators, the following platforms were used:

- **Platform 1:** One of the nodes of the cluster Guane of Supercomputing Center of Universidad Industrial de Santander with the following configuration: two AMD EPYC 9554 64-Core (two threads per core) @ 3.1 GHz Processors and 375 GB of RAM memory.
- **Platform 2:** A workstation with One Intel(R) Xeon(R) E-2136 CPU 6-Core (two threads per core), @ 3.30GHz processor with 32 GB of RAM, and a NVIDIA Corporation GP106GL Quadro P2000 5GB. This node is used only for GPU-capable simulators.

### 3.3. Intel-QS

It is an open-source quantum circuit simulator implemented in C++. It uses multi-processing and has an intuitive Python interface. It is a full-state vector simulator using

arbitrary single-qubit gates and gates controlled by two qubits [17]. The Intel Quantum Simulator leverages the full capabilities of an HPC system through its shared and distributed memory implementation. The implementation on a single node incorporates enhancements such as vectorization, threading, and cache optimization through the process of gate fusion. The primary object in the Intel Quantum Simulator (IQS) is the QubitRegister, representing the quantum state of the qubits in the system of interest. When declaring a QubitRegister, the number of qubits must be specified to allocate enough memory to describe their state. The state can then be initialized to any computational basis state, uniquely identified by its index. In this work, we used the C++ interface of this simulator.

### 3.4. Quantum++

It is a general-purpose multithreaded quantum simulator with high performance. The library is not restricted to qubit systems or specific quantum information processing tasks, being capable of simulating arbitrary quantum processes [18]. Quantum++ is developed using standard C++17 and has minimal external dependencies. It primarily utilizes the Eigen 3 linear algebra template library, which is header-only. Additionally, when available, it employs the OpenMP library to facilitate multiprocessing. The primary data types are complex vectors and complex matrices, such as complex dynamic matrices, double-precision dynamic matrices, complex dynamic column vectors, complex dynamic row vectors, etc.

### 3.5. qsim

Developed by Google, qsim is an optimized quantum circuit simulator that uses gate fusion and vectorized instructions to simulate up to 40 qubits on a powerful workstation [19]. Integrated with Cirq, it provides a robust environment for developing and testing quantum algorithms. To achieve cutting-edge simulations of quantum circuits, it uses gate fusion, AVX/FMA vectorized instructions, and openMP multithreading. This relies on cuQuantum to integrate GPU support.

### 3.6. cuQuantum

NVIDIA's cuQuantum SDK is another leading tool, designed to accelerate quantum circuit simulations on GPUs. This toolkit is essential for developers looking to leverage the power of GPUs to enhance simulation performance and scalability. It provides an integrated programming model tailored for a hybrid environment, enabling the combined operation of CPUs, GPUs, and QPUs.

### 3.7. QuEST

QuEST, or the Quantum Exact Simulation Toolkit, is a high-performance open-source quantum computing simulator designed for simulating quantum circuits, state-vectors, and density matrices. Developed by the Quantum Technology Theory Group at the University of Oxford, QuEST is distinguished by its ability to utilize multithreading, GPU acceleration, and distribution, making it highly effective across various computing environments, from laptops to networked supercomputers. The toolkit is capable of simulating both pure quantum states and mixed states with precision, and supports a wide array of quantum operations. It allows simulations that are extensible and adaptable, thanks to its open-source nature and support for various back-end hardware via its simple and flexible interface [20]. QuEST represents a pure state for a system of  $n$  qubits using  $2^n$  complex floating-point numbers, with each real and imaginary component having double precision by default. However, QuEST can be configured to use single or quad precision if desired. The simulator stores the state using C/C++ primitives, which means that by

default, the state vector alone consumes  $16 \times 2^n$  bytes of memory, where  $n$  stands for the number of qubits.

### 3.8. Qrack

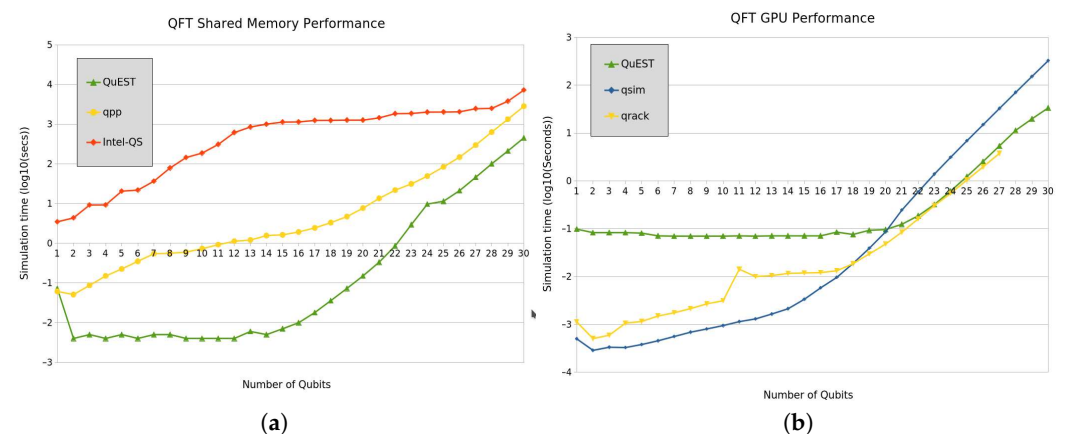
Qrack is a high-performance quantum computer simulator that is written in C++ and supports OpenCL and CUDA [21,22]. It is particularly notable for its ability to simulate arbitrary numbers of entangled qubits, limited only by system resources. Qrack is designed to be embedded in other projects and includes a comprehensive suite of standard quantum gates, along with variations suitable for register operations and arbitrary rotations. The simulator is integrated with other quantum computing frameworks like ProjectQ and Qiskit, enhancing its versatility and application. Qrack also features optimizations for noiseless pure state simulations and includes tools that aid in the control, extension, and visualization of data from quantum circuits. Qrack maintains the state representation in a factorized form to enhance simulation efficiency. A general ket state  $|\psi\rangle$  of  $n$  qubits is described by  $O(2^n)$  complex amplitudes.

### 3.9. Simulator Comparison

Regarding academic, community, and industry support for these simulators, the continual updates, active support, and documentation for these tools contribute to their status as leaders in the field. This ongoing development ensures they remain relevant and valuable as quantum computing technology evolves. Each of these simulators offers unique features and optimizations, making them suitable for various aspects of quantum computing research and development. Their continual evolution is critical as the quantum computing field strives to solve more complex problems and improve algorithm efficiency. Table 1 shows a comparison of the evaluated simulators of their design properties and optimization mechanisms.

Other projects, like XACC and Qiskit, provide a full-stack approach to quantum computing, including simulators, compilers, and the possibility to run programs on real quantum processors.

For convenience and agility, those simulators that provided QFT in their examples were compared under equal conditions. First, OpenMP capable simulators are shown in Figure 3a. Second, GPU capable simulators are depicted in Figure 3b.



**Figure 3.** Comparison of the quantum Fourier transform using different simulators and optimization techniques. (a) Shared Memory Performance using platform 1. (b) GPU Performance using platform 2.

**Table 1.** Leading open-source simulators for quantum computing comparison.

Features	Intel-QS	Quantum++	QuEST	Qrack	qsim
Optimization Techniques	Uses MPI for distributed computing, optimizes for multi-core and multi-nodes, supports MKL for mathematical operations	Employs template metaprogramming for compile-time optimizations, supports OpenMP for parallelization	Utilizes GPU acceleration, multi-threading, and can be distributed across networked supercomputers	Leverages gate fusion and OpenCL for parallel execution across different hardware platforms	Optimizes simulations with gate fusion, AVX/FMA vectorized instructions, and OpenMP
Hardware Support	Optimized for high-performance computing systems, can be deployed on cloud infrastructures	Compatible with various architectures via C++ standardization, no specific hardware acceleration mentioned	Supports laptops to supercomputers, compatible with GPUs and distributed systems	Designed for broad compatibility with OpenCL-supporting GPUs and CPUs	Runs efficiently on high-core-count CPUs and potentially on any system supported by Cirq
Programming Model	Provides C++ and Python interfaces, supports state vector simulation	C++ library with emphasis on flexibility and ease of integration	Offers a C library that is easy to integrate and extend, with optional Python bindings	C++ based, with a focus on integrating with other quantum computing frameworks like Qiskit	Integrated with Cirq, emphasizes ease of use in Python for simulating quantum circuits
Design Properties	Focuses on scalability and performance across different computational environments	Prioritizes modular, generic programming for ease of adaptation and maintenance	Designed for precision and versatility in quantum state manipulation	Prioritizes rapid prototyping and flexibility for embedding in various applications	Designed to simulate large quantum circuits with high precision
Unique Features	Supports dynamic circuit simulation and state manipulation during runtime	Highly adaptable to various quantum computing models due to generic programming approach	Extensible and supports detailed state analysis tools like fidelity and entanglement measures	Integrates classical computing elements within quantum simulations for enhanced functionality	Deeply integrated with Google's quantum computing framework, providing extensive simulation capabilities

#### 4. Implementing a Simulator

To gain a deeper understanding of the fundamental operations of quantum computing and to test the various memory management approaches, a software quantum simulator prototype was developed in C++ (The Memory efficient Quantum Simulator, TMFQS) [23]. This prototype was designed in such a way that it allows us to change strategies easily through minimal modifications. It allows us to easily adjust the data structures to represent the fundamental concepts of quantum computing and the use of compression libraries. On the other hand, to demonstrate the construction of a software quantum simulator in a simple way, the scope of this work was limited to optimization techniques using shared memory and distributed memory. The implementation using GPUs is left for further work.

It has to be pointed out that this prototype does not implement all the concepts of quantum computing, such as quantum error correction, entanglement, measurement, and an extended set of quantum gates. The measurement operation, a fundamental aspect of quantum computing, was not implemented in this prototype because the primary focus of this research was to evaluate and optimize memory management strategies for quantum state simulation. The objective was to explore various methods, such as state pruning, data compression, and parallelization, to enhance the efficiency of memory use in large-scale quantum simulations. Since these techniques do not inherently depend on measurement processes, implementing a measurement function was not essential to achieving the research goals. However, the measurement operation could be incorporated in future iterations to extend the simulator's capabilities for practical quantum algorithm execution. Several scenarios were implemented to carry out the tests.

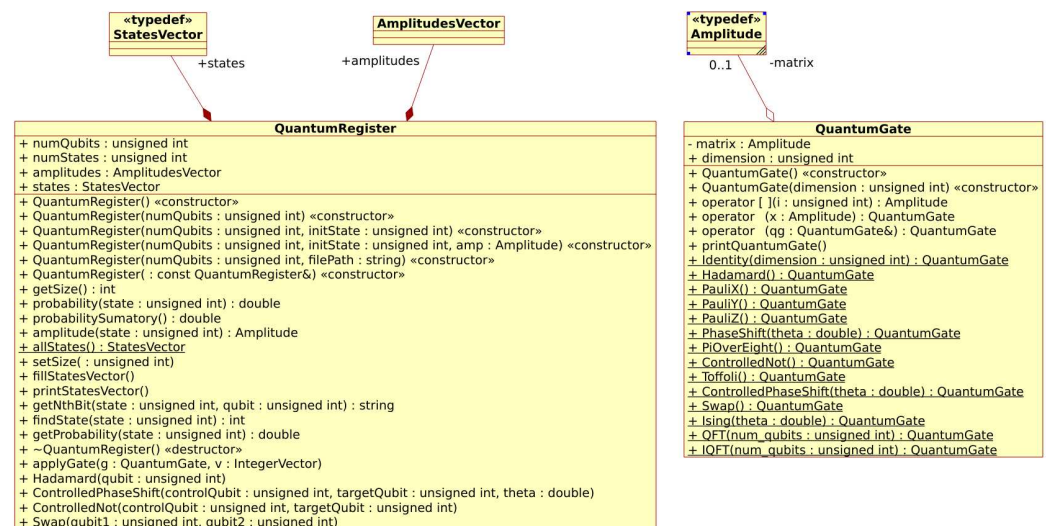
- Dynamic memory management: The primary purpose is to test the strategy of removing less probable states.
- Full State: The objective is to accelerate the simulations, avoiding the overhead introduced by the search for the states.
- Full State with OpenMP: The intention is to accelerate the simulations of the previous version.
- Full State with data compression: The purpose is to test a lossy compression library like Compressed Floating-Point library (ZFP).
- Full State with MPI: The main objective of this scenario is to distribute the amplitude vector among different computing nodes, allowing for a greater number of qubits.
- Full State with MPI and data compression: Here, data compression was incorporated into the previous scenario.

#### 4.1. Representation of Fundamental Concepts

As we saw in the previous section, the basic simulation concepts include the following elements.

- Quantum Register: Comprised of a state vector and an amplitude vector.
- Quantum Gates: Matrix representation of quantum gates. Only one-qubit and two-qubit gates were considered.
- Applying quantum gates to a quantum register.

We have used an array of double-precision floating-point numbers to store the amplitudes. No data structure has been used to represent the states, since the vector indices are used to refer to them. To implement the method to apply a quantum gate to a quantum register (`QuantumRegister::applyGate`), we use the technique represented in Equation (5) proposed by [24]. In Figure 4, we can observe the main classes of the prototype.



**Figure 4.** Class diagram of the prototype: The `QuantumRegister` class represents a quantum state and implements the main method to transform a quantum state (`applyGate`). The `QuantumGate` class implements a small set of quantum gates using the matrix representation.

#### 4.2. Single Processor Case

To evaluate the quantum simulator's performance and memory management strategies, we begin by analyzing its execution on a single processor. This section focuses on the core elements of such simulations, particularly the representation of the state vector and the application of quantum gates.

#### 4.2.1. State Vector

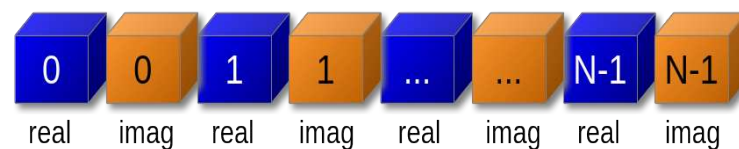
In this subsection, we discuss how the simulator represents the quantum state. Specifically, we describe the structure of the state vector, which stores the amplitudes of all possible quantum states, and explain how these amplitudes are organized to optimize memory usage and computational efficiency.

The state vector is a linear combination of states represented by the following expression.

$$|\psi\rangle = \alpha_0|0\dots00\rangle + \alpha_1|0\dots01\rangle + \dots + \alpha_{2^n-1}|1\dots11\rangle \quad (2)$$

where  $\alpha_i$  are the amplitudes. As we said previously, these amplitudes are complex numbers, so we need two `float` or two `double` numbers to represent them in the code. Of course, the state vector must fit in the local memory.

The amplitudes of the states are implemented using a single-dimensional double-precision array stored in a continuous memory space. To increase performance, a single array was used to store both the real and the imaginary parts of each amplitude; that is, the state vector was linearized. The real parts are placed in the odd positions of this arrangement, and the imaginary parts are placed in the even positions. This strategy avoids jumping between two arrays, one for the real part and one for the imaginary part. Figure 5 depicts this data structure.



**Figure 5.** Linearized state vector.

#### 4.2.2. Quantum Gates

Like other simulators such as Intel-QS, this prototype only implements single-qubit gates and controlled two-qubit gates. The minimum list of quantum gates developed to implement the quantum Fourier transform algorithm is as follows: Identity, Hadamard, ControlledPhaseShift, ControlledNot, Swap. All these quantum gates are implemented as two-dimensional double-precision arrays. This reduced set of quantum gates limits the simulation of algorithms that require additional gates. However, adding either new single-qubit or controlled two-qubit gates is very easy. Just insert the corresponding matrix into the `quantumGate.cpp` source file.

#### 4.2.3. Applying a Quantum Gate to a Quantum Register

To apply a quantum gate  $G_k$  to the  $k$ -th qubit of a state vector  $|\psi\rangle$ , we have the following result.

$$G_k|\psi\rangle = |\psi'\rangle = \begin{pmatrix} \alpha'_{0\dots00} \\ \alpha'_{0\dots01} \\ \vdots \\ \alpha'_{1\dots10} \\ \alpha'_{1\dots11} \end{pmatrix} \quad (3)$$

#### Applying Single-Qubit GATES

The first traditional approach to face this problem is using sparse matrix management methods. However, [24,25] state that applying a single-qubit gate  $G_k$

$$G_k = \begin{pmatrix} g_{00} & g_{01} \\ g_{10} & g_{11} \end{pmatrix} \quad (4)$$

to the  $k$ -th qubit of a quantum register of  $N$  qubits is equivalent to applying the gate to pairs of amplitudes whose indices differ by  $k$ -th bits from their binary index.

$$\begin{aligned} \alpha'_{*...*0_k*...*} &= g_{11} \cdot \alpha_{*...*0_k*...*} + g_{12} \cdot \alpha_{*...*1_k*...*} \\ \alpha'_{*...*1_k*...*} &= g_{21} \cdot \alpha_{*...*0_k*...*} + g_{22} \cdot \alpha_{*...*1_k*...*} \end{aligned} \quad (5)$$

This implies that all state vector elements must be processed in pairs. Let us see how to apply the Hadamard gate to the first qubit of the state  $|00\rangle$ . For the values  $k = 0$ ,  $*...*0_k*...* = 00$ ,  $*...*1_k*...* = 10$ ,  $\alpha_{00} = 1 + 0i = 1$ ,  $\alpha_{10} = 0$ , and Hadamard gate,

$$H = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad (6)$$

Replacing these values in Equation (5), we obtain the following results.

$$\begin{aligned} \alpha'_{00} &= \frac{1}{\sqrt{2}} \cdot 1 + \frac{1}{\sqrt{2}} \cdot 0 = \frac{1}{\sqrt{2}} \\ \alpha'_{10} &= \frac{1}{\sqrt{2}} \cdot 1 - \frac{1}{\sqrt{2}} \cdot 0 = \frac{1}{\sqrt{2}} \end{aligned} \quad (7)$$

#### Applying Two-Qubits Gates

Similarly, to apply a controlled two-qubit quantum gate to a quantum register, using a control qubit  $c$  on a target qubit  $t$ , authors of [25] state that the new amplitudes can be obtained by performing the following operations:

$$\begin{aligned} \alpha'_{*...*1_c*...*0_t*...*} &= g_{11} \cdot \alpha_{*...*1_c*...*0_t*...*} + g_{12} \cdot \alpha_{*...*1_c*...*1_t*...*} \\ \alpha'_{*...*1_c*...*1_t*...*} &= g_{21} \cdot \alpha_{*...*1_c*...*0_t*...*} + g_{22} \cdot \alpha_{*...*1_c*...*1_t*...*} \end{aligned} \quad (8)$$

Let us see how to apply the Controlled Phase Shift (CPS) gate to the second qubit of the state  $|11\rangle$  controlled by the first qubit. All amplitudes are equal to 0 except  $\alpha_{11}$ , which is equal to 1. Replacing these values in the Equation (8), we have

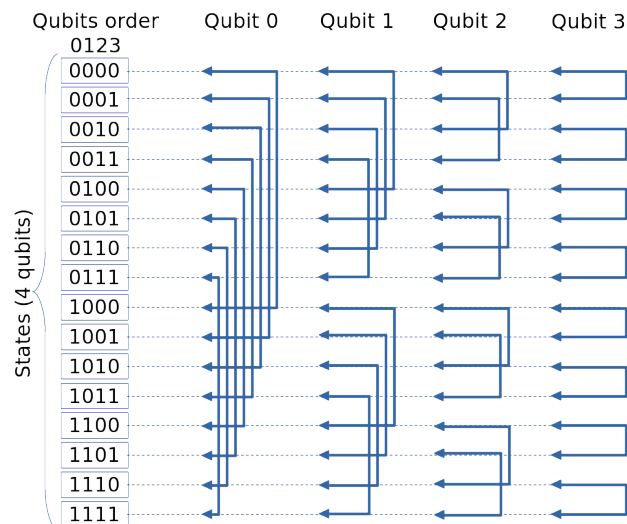
$$\begin{aligned} \alpha'_{10} &= 1 \cdot 0 + 0 \cdot 0 = 0 \\ \alpha'_{11} &= 0 \cdot 1 + e^{i\phi} \cdot 1 = e^{i\phi} \end{aligned} \quad (9)$$

Thus, we obtain the amplitude values for the states  $|10\rangle$  and  $|11\rangle$ .

#### Qubits Order

Some simulators, like qiskit, reverse the order of the qubits such that qubit 0 corresponds to the least significant bit of the binary representation of the state. In this case, the distance between  $\alpha'_{*...*0_k*...*}$  and  $\alpha'_{*...*1_k*...*}$  is equal to  $2^k$ .

In this work, we maintain the natural order of the qubits. For example, in state  $|011\rangle$ , qubit 0 is the leftmost. Therefore, the distance between  $\alpha'_{*...*0_k*...*}$  and  $\alpha'_{*...*1_k*...*}$  is equal to  $2^{(numQubits-1)-(k-th\ qubit)}$ . To illustrate this, Figure 6 shows the distance between the states of a four-qubit state vector.



**Figure 6.** State pairs affected by qubit operation.

Generally, a single-qubit gate can be applied to a quantum register, performing the pseudo-code showed in Listing 1.

**Listing 1.** Applying a single-qubit gate.

---

```

For each amplitude in the state vector
do
    Calculate the new amplitudes for the current state.

    Find the pair corresponding to the current quantum state

    Calculate new values for the amplitude of the new state.

done

```

---

In summary, calculating the amplitudes for the current state and the new affected state is performed as follows: Determine the value of the current state's amplitude using Equation (5). Then, find the pair corresponding to the current state, and finally, calculate the value of the latter using that same equation.

To find the pair corresponding to the current state, we can use two methods: the first calculates the distance using the relation  $2^{(numQubits-1)-(k-th\ qubit)}$ , as we explained before. The second method applies an XOR operation between the binary representation of the current state and a sequence of 0s with a 1 placed in the k-th position corresponding to the qubit we are working on. For example, applying a quantum gate on the 0th qubit on a four-qubit state  $|0101\rangle$ , we can find the corresponding pair using the following operation.

$$\begin{array}{r}
 0101 \\
 \oplus 1000 \\
 \hline
 1101
 \end{array} \tag{10}$$

This result can be corroborated in Figure 6. C++ offers binary operations to execute this operation efficiently and effortlessly. This can be observed in Listing 2.

**Listing 2.** C++ XOR operation.

---

```

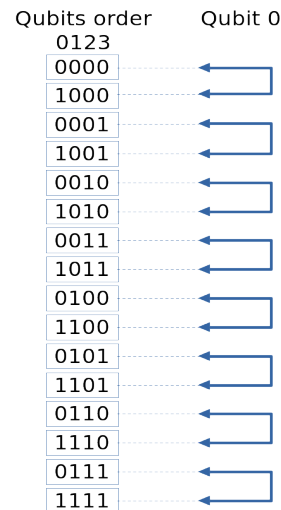
unsigned int pos = numQubits - qubit - 1 ;
unsigned pairState = currentState ^ ((uint)1 << pos);

```

---

#### 4.2.4. State Pruning

In the version of the simulator where the least probable states are pruned, we use dynamic memory management because the states are stored non-sequentially in memory. This arrangement results from their computation via Equation (5). Consequently, a state search method was developed to facilitate access to a specific state for calculations in subsequent iterations. However, performance is negatively impacted because a lot of time has to be spent searching for a state's values before they are used in a calculation. Figure 7 depicts the order of a three-qubit state vector after applying a single-qubit gate (Hadamard) on qubit 0.



**Figure 7.** State vector order using dynamic memory approach.

Because of this, the less probable states elimination approach was discarded early; therefore, we focus on pure states, which implies that the state vector contains the complete information about the quantum state, and this approach was adopted for the rest of this simulator's design.

#### 4.2.5. Compressing the State Vector

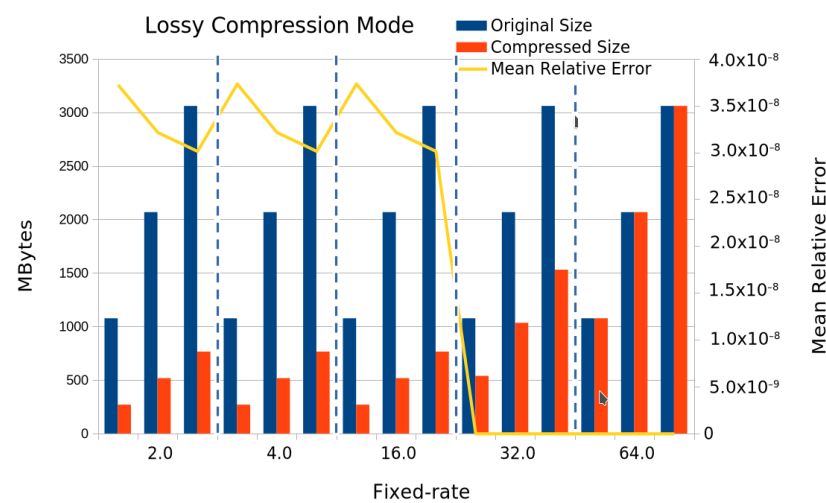
The large volumes of data produced by extreme-scale scientific research and applications have driven the development of various data compression techniques for years. The compression methods are optimized for floating-point data. However, they require additional calculations to compress and decompress data before working with it. Leading data compressors can be classified into two categories:

- **Lossless Compressors:** these compressors preserve all the data.
  - They use variable-length encoding algorithms such as Huffman encoding, Arithmetic encoding, and Dictionary encoders.
  - They often produce low compression ratios, typically around 2:1.
- **Error-Bounded Lossy Compressors:** allow some controlled distortion. They can be broadly classified into
  - The data-prediction-based compression model;
  - The domain-transform-based compression model.

A key objective of this work is to identify a method for reducing memory consumption in a software quantum simulator. To achieve this, we have chosen error-bounded lossy compressors, the compression technique that offers the best compression rate.

To compress the amplitude vector, we use the ZFP library [26] as it provides significant performance in accuracy and data size reduction. Although ZFP supports both lossy and lossless compression, as we stated before, we used the lossy approach to gain a better compression rate. To support this claim, both ZFP compression algorithms were tested.

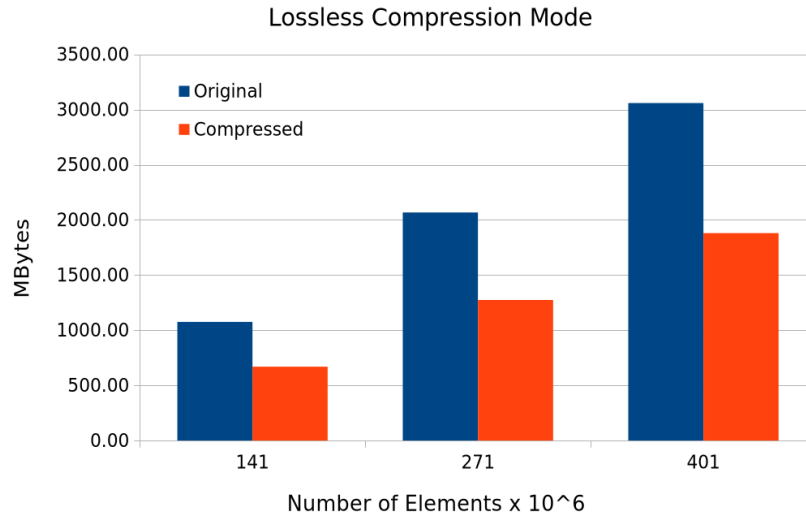
The graph in Figure 8 illustrates the compression of vectors of sizes 1, 2, and 3 GB at various compression rates using the fixed-rate mode. At fixed rates of 2.0–16.0, compression is strong and stable. At 32.0, the size doubles, and at 64.0, the size is almost the same as the original (with no actual compression). Lower fixed rates (2–16) yield significant savings (over 75%), while high fixed rates (32–64) sharply reduce compression efficiency. The highest error occurs at fixed rates 2.0–16.0, at 32.0, the error drops drastically, and at 64.0, the error is zero. The chart effectively visualizes a compression–quality trade-off: compression error decreases as the fixed rate increases. For lower fixed rates, we obtain better compression but worse accuracy, and for higher fixed rates, we obtain worse compression but better accuracy.



**Figure 8.** Lossy compression mode.

The graph in Figure 9 shows the compression of an array using the lossless mode of ZFP. We can observe that the compression ratio is very consistent across different vector sizes (1.61–1.63×), indicating that ZFP’s lossless mode behaves predictably regardless of vector length, at least in this small sample. The relationship between the number of elements and file size is roughly linear in both original and compressed sizes, confirming that ZFP’s lossless encoding scales proportionally with input data. The average compression ratio of 1.62x is modest. This is expected for lossless compression of floating-point data. If the algorithm is sensitive to precision and cannot afford lossy compression, this trade-off (38% space saving) may still be worthwhile, especially in large-scale simulations or scientific data storage.

To go from a vector of amplitudes using traditional data types to a compressed vector, change the corresponding line in the types.h source file from `typedef std::vector<double> AmplitudesVector;` to `typedef zfp::array1<double> AmplitudesVector;` of course, the corresponding header file from the ZFP library must be included.



**Figure 9.** Lossless compression mode.

#### 4.2.6. Shared Memory Case

To improve performance, parallelizing the code is necessary. The first method is to apply a shared memory programming model. This was achieved using OpenMP.

We use `valgrind` to run program profiling and determine the sections of code that consume the most resources. Afterward, it was determined that the `QuantumRegister::applyGate` method is the component of the simulator where we had to focus on increasing performance. Figure 10 shows the profiling results.

Incl.	Self	Called	Function
100.00	0.00	(0)	0x00000000000020ed0
92.77	0.00	1	(below main)
92.77	0.00	1	__libc_start_main@@GLIBC_2.34
92.77	0.00	1	(below main)
92.77	0.00	1	main
78.24	0.00	1	quatumFourierTransform(QuantumRegister*)
78.13	21.55	96	QuantumRegister::applyGate(QuantumGate, std::vector<uns
64.67	0.00	6	QuantumRegister::Swap(unsigned int, unsigned int)
64.66	0.00	18	QuantumRegister::ControlledNot(unsigned int, unsigned int)
25.83	17.50	1 589 405	mcount
16.20	5.28	671 532	QuantumGate::operator[](unsigned int)
15.85	7.79	495 459	copyBits(int, int, int, int)

**Figure 10.** Simulator profiling.

The `QuantumRegister::applyGate` method iterates through the state vector, implementing Equation (5). To enhance performance, we partition the data and execute instructions on segments of the state vector, thereby speeding up the simulation. It is crucial to carefully manage the method's internal variables to prevent race conditions.

#### 4.3. Distributed Memory Case

In the distributed memory model, the state vector needs to be divided among  $numProcs$  processes. On the other hand, Equation (5), proposed in [24], indicates that the calculation of the amplitudes of the states must be performed in pairs; therefore, we must guarantee that the number of amplitudes per process is even. To achieve this, we use the relationship

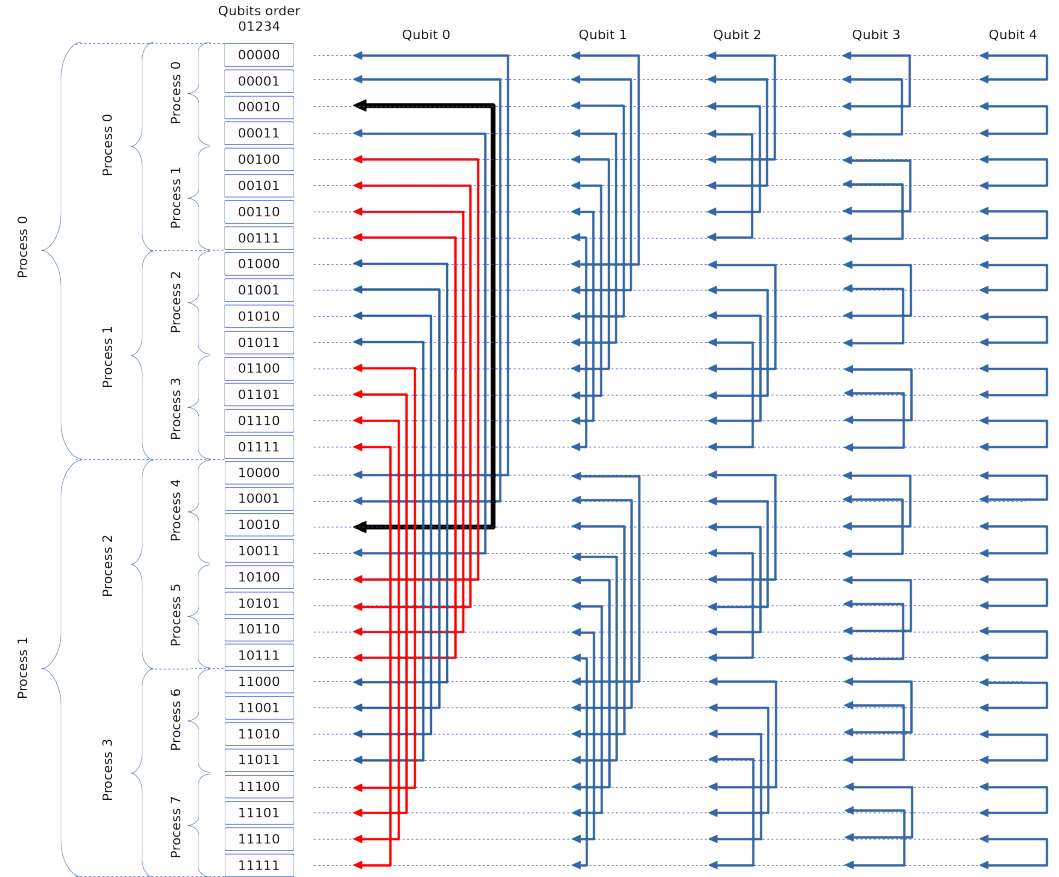
$$numProcs = \frac{2^{numQubits}}{2^m} \quad (11)$$

where  $2^m$  is the number of states per process. In this case, we can face two cases:

- The pair corresponding to the current state is locally stored.

- The pair corresponding to the current state is located in another process. In this case, it is necessary to communicate values and results.

Figure 11 shows the pairwise calculation scheme for a five-qubit state vector, applying each qubit. Partitioning with two, four, and eight processes is also shown to visualize the communication process easily when we use the distributed programming model.



**Figure 11.** Pairwise calculation scheme for a 5-qubit state vector.

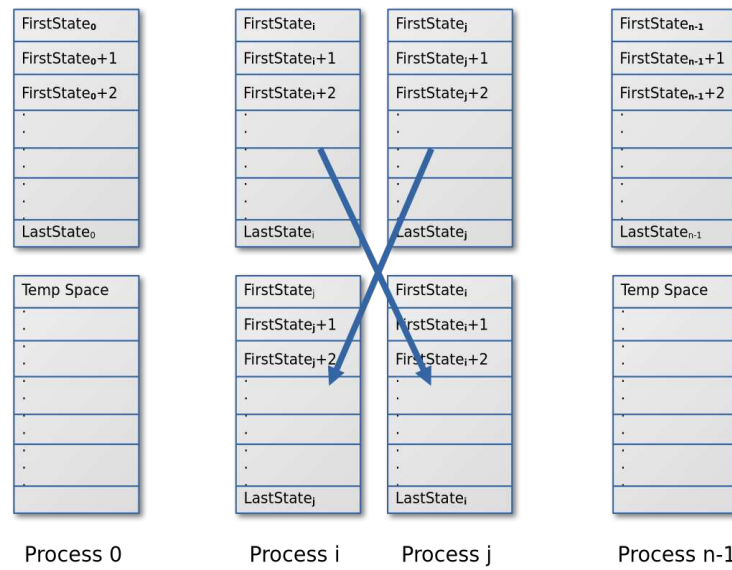
For instance, consider performing a calculation on qubit 0 of the state  $|00010\rangle$ ; the corresponding pair would be  $|10010\rangle$ . If two processes are used, communication should be established with process 1. If four processes are utilized, the remote process is process 2. Lastly, if eight processes are employed, the remote process will be process 4.

We use the following expression to calculate the process’s identifier, where the corresponding pair is located.

$$remoteProcID = \frac{pairState}{2^m} \tag{12}$$

In Figure 11, it is evident that for two processes, specifically regarding qubit 0, the number of communications required is  $2^{\text{numQubits}/2}$ . This substantially degrades performance. To mitigate the overhead caused by the extensive number of communications, the entire segment of the state vector is exchanged between the peer processes involved, as outlined in Equation (5). The calculations are then made locally, and the results are communicated back to the original process.

For this reason, we are unable to use the total sum of local memory of each node to augment the number of qubits, and can only utilize half of the combined memory from all nodes. Figure 12 depicts this idea.



**Figure 12.** Data exchange between processes.

Combining amplitude vector compression with amplitude vector distribution across multiple processes is an approach that can be effective both in terms of efficient memory usage and overall simulation performance. The version where a compressed vector is used to store the amplitudes was parallelized to achieve this. To obtain effective performance, the state vector portions are transmitted in a compressed manner. This makes communications faster because the message size is reduced.

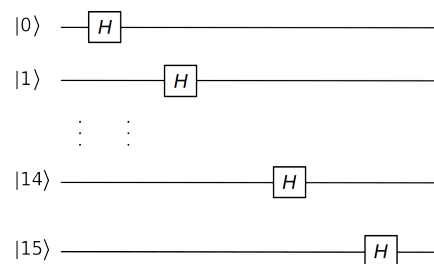
To achieve compressed messaging, the compressed portions of the state vector must be serialized, and a custom MPI data type must be used in send and receive functions.

#### 4.4. Simulator Verification

In order to validate the accuracy of our quantum simulator, we have executed different tests and compared the outputs with intel-qc and quantum++.

#### Quantum Gates Tests

To test the superposition principle, we apply a Hadamard gate to a quantum register of four qubits. The quantum circuit in the Figure 13 was used.



**Figure 13.** Test quantum circuit.

The test was executed by initializing the first state with a probability equal to one, that is to say,  $1 \times |0000\rangle$ . Then, we repeat the experience with  $1 \times |0001\rangle$  and so on until executing the test with the last state  $1 \times |1111\rangle$ . The results of executing this quantum circuit with Intel-qc, quantum++, and TMFQS were the same.

## 5. Results

This section presents the results of several quantum simulation tests performed using the TMFQS software quantum simulator developed in C++. By simulating fundamental quantum operations, we can assess how well these strategies reduce memory consumption and improve the efficiency of quantum computing simulations on classical hardware. Throughout the section, we compare the simulator's performance with and without the proposed memory management techniques, highlighting the improvements achieved. By providing a comprehensive evaluation of these memory management strategies, this section aims to contribute to ongoing efforts to make quantum computing simulations more efficient and scalable, ultimately advancing the field of quantum computing.

### 5.1. Algorithm Selected for Testing

TMFQS was evaluated using the quantum Fourier transform, as in assessing the simulators presented previously.

### 5.2. Test Platform

To run the simulations, we use two high-performance nodes from the scientific computing center of the Universidad Industrial de Santander (SC3-UIS) with the following characteristics: two AMD EPYC 9554 64-Core (two threads per core) @ 3.1 GHz Processors and 375 GB of RAM.

#### 5.2.1. States Pruning

We can free up memory that is not needed by eliminating the quantum states that have the smallest chance of occurring; that is, eliminating those states whose amplitude is close to zero. However, in addition to all the points against this approach, like loss of fidelity, impact on algorithm accuracy, error accumulation, threshold sensitivity, and impact on quantum entanglement, this requires dynamic memory management, which introduces a lot of extra work because before applying a quantum gate to a state, we need to search for it in the array because the states are not ordered. That is to say, to apply every quantum gate, we need to execute  $2^n$  search operations. This approach was tested using the quantum Fourier transform algorithm. The quantum register contains all the states at the end of executing the quantum Fourier transform algorithm. Due to the initial superposition process, the quantum register also has all the states in the first stage of Grover's algorithm. Therefore, this approach does not work well for these algorithms.

For these reasons, along with the risks outlined previously, we have decided to discard this approach because the disadvantages significantly outweigh the benefits.

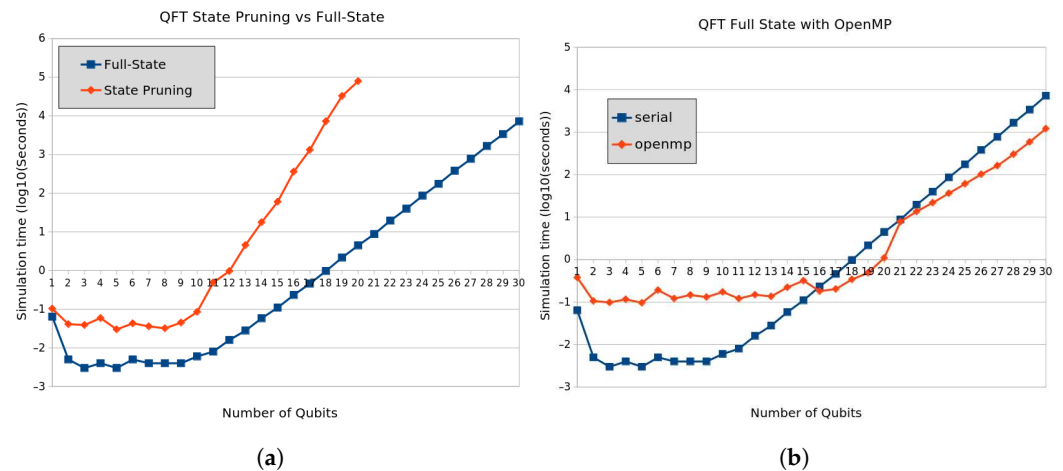
#### 5.2.2. Full-State Quantum Register

A quantum register with all the states arranged in a sequence can reduce the overhead of searching for quantum states. This also eliminates the need for an extra data structure to store the states and uses the indices of the amplitude vector to handle the quantum states. The total amount needed using this strategy is  $2^{numQubits} * 16$  bytes for double precision floating-point numbers. However, we save  $2^{numQubits} * 4$  bytes, avoiding the state vector array.

The graph of Figure 14a shows the performance of QFT applying state pruning (dynamic memory) vs full-state strategies. Simulations using dynamic memory involving more than 20 qubits were discarded due to their execution time exceeding one day.

Exponential growth is observed from an 18-qubit state vector in the dynamic memory approach. This is a consequence of a substantial increase in the memory needed to represent the state in question and, therefore, the processing time required. As can be seen by

comparing these strategies, the workload overhead using dynamic memory is significant. We have parallelized the full-state version to increase performance using the shared memory model with OpenMP. In the graph of the Figure 14b, we can see the results.



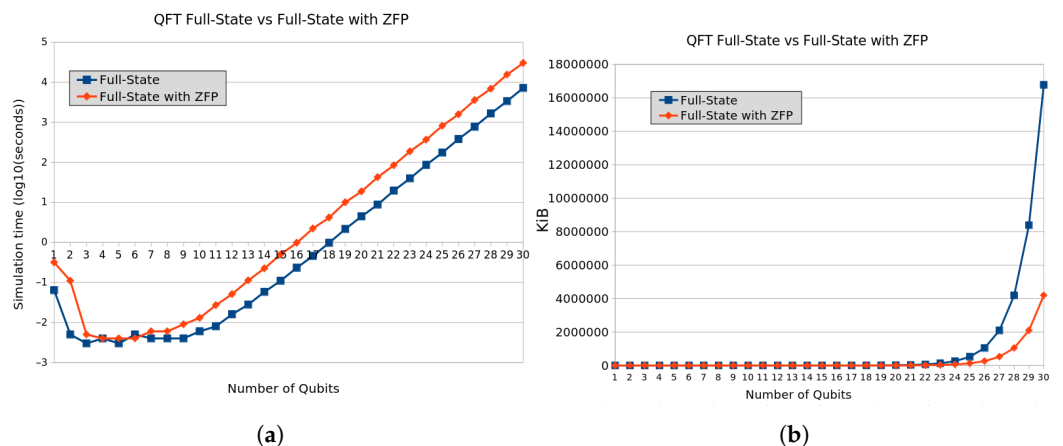
**Figure 14.** QFT performance with different approaches. (a) Dynamic Memory vs. Full-State Approach. (b) QFT Full-State Approach with OpenMP ( $\log_{10}$ ).

We observe a significant decrease in the processing time between the serial and parallel execution of the full-state version as the number of qubits increases. That is, a considerable acceleration is obtained by parallelizing the simulation. However, for a clearer interpretation of the results, we show the results calculating the base 10 logarithm of the simulation time. In the graph shown in Figure 14b, it is evident that for smaller numbers of qubits, there is an overhead caused by the setup of the parallel environment.

### 5.2.3. Data Compression

We have selected one of the most widely used C++ libraries for data compression, ZFP, to test this approach. We modified the full-state version of the simulator to compress the amplitude vector. The graph of Figure 15a shows the performance comparison between full-state vs full-state using ZFP. The base-10 logarithm is used to more clearly highlight the differences between the two simulations. It is evident that the overhead introduced by the compression and decompression process is substantial.

The graph of Figure 15b shows the amount of memory used by both simulator versions. We can observe that the compression approach is highly efficient. This enables the possibility of increasing the number of qubits in simulations.

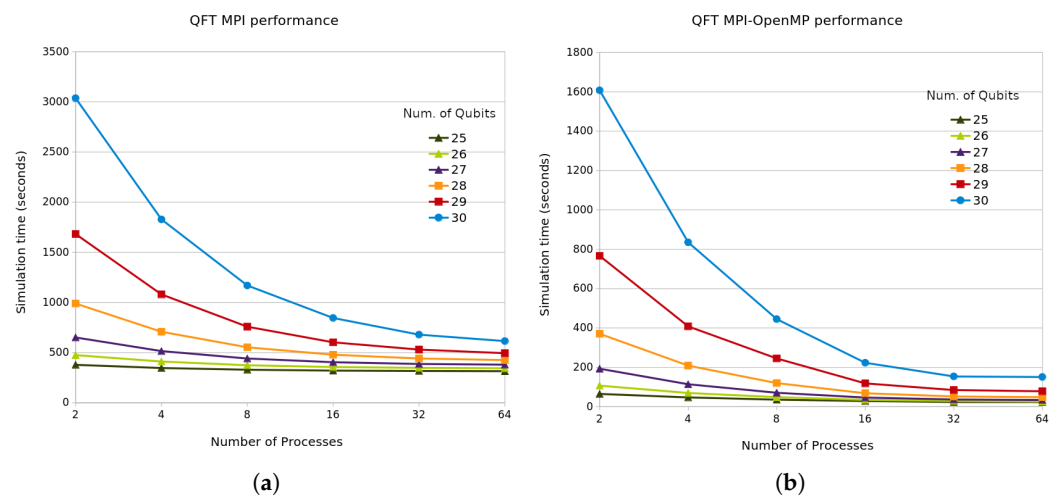


**Figure 15.** QFT performance with ZFP. (a) QFT Full-State with ZFP ( $\log_{10}$ ). (b) QFT Full-State with ZFP Data Size.

#### 5.2.4. Distribute the Quantum Register Across Multiple Nodes

We have developed a simulation version employing MPI to increase memory capacity by leveraging the RAM of additional computer nodes. To uphold computational efficiency, it is essential to underscore the necessity of maintaining an optimal ratio between the number of processes and the allocation of qubits per process. The graph shown in Figure 16a demonstrates the performance of the quantum Fourier transform across a range of qubit counts from 25 to 30, using 2, 4, 8, 16, 32, and 64 processes. The relationship of Equation (11) is valid only from 7 qubits onwards; however, to visualize the performance more clearly, we use the range of 25 to 30 qubits. In this case, a logarithmic scale was unnecessary since the graphs were clear enough.

In addition to achieving better performance, we can see that by increasing the number of processes, we can increase the number of qubits and reduce the size of messages required to exchange partial results between processes. We can see also that parallelism is helpful for a number greater than 26 qubits.



**Figure 16.** QFT performance using parallel techniques. Each line corresponds to a specific number of qubits. (a) QFT Full-State with MPI. (b) QFT Full-State with MPI and OpenMP.

#### 5.2.5. Combination of Distributed Memory and Shared Memory Approaches

We have developed a simulator version that combines MPI with OpenMP to achieve better performance. In this approach, the state vector is evenly distributed across the processes using MPI. OpenMP is then employed to parallelize the applyGate method, further enhancing performance. The graph in Figure 16b illustrates the performance of the quantum Fourier transform using this hybrid approach. Once again, the graph is clear enough; therefore, a logarithmic scale is not necessary.

Comparing the results of the graphs in Figure 16a,b, we see that the combination of MPI and OpenMP increases the performance, especially for cases where the size of the state vector portion at each node is large.

#### 5.2.6. Combination of Distributed Memory and Data Compression Approaches

Taking advantage of distributed resources to have more memory available, combined with data compression, it makes it possible to perform simulations with a larger number of qubits.

We have already seen that there is a processing overhead introduced by the compression process, however, the transmission of compressed data contributes positively to overall performance. We have modified TMFQS to test this approach. Figure 17 shows the execution for a range of 25 to 30 qubits with a variation in the number of processes equal

to 2, 4, 8, 16, 32, 64. It shows the performance of the distributed memory approach with data compression.

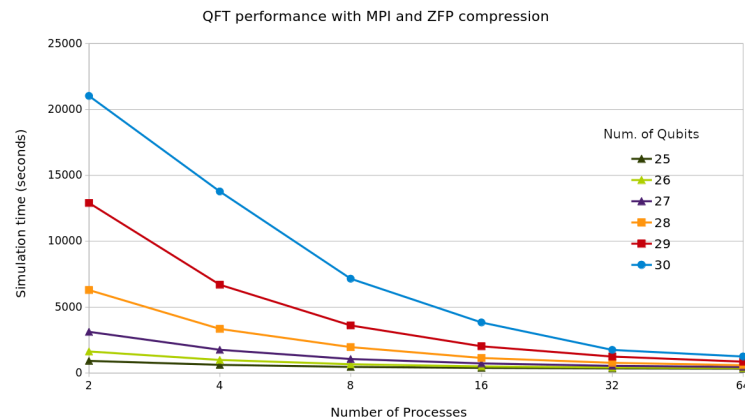


Figure 17. QFT with MPI with data compression.

From the graph in Figure 17, we can see that performance has decreased; however, the reduction of the required memory is significant because the same strategy used in the section on data compression is adopted here. It has to be pointed out that this strategy is valid only if portions of the quantum register are transmitted in a compressed form.

### 5.2.7. Quantum Simulators Comparison

To validate the results obtained with TMFQS, a comparison is made with other simulators. First, specific conditions must be established to allow a fair comparison of the simulators studied. The common conditions were using a single computing node with a shared memory model. Simulators that use GPUs are excluded because their performance is much higher than the others, but their scaling is limited. The case of using distributed memory is also excluded because only some include this capability. The graph in Figure 18 shows the performance of the quantum Fourier transform for intel-qs, quantum++, QuEST, and TMFQS using the shared memory model. All simulators were compiled with the Intel OneAPI suite and traditional optimization flag (-O2). The selected simulators use the C++ complex numbers data type for memory management. They use a full-state vector scheme, which allows better performance but does not reduce memory consumption.

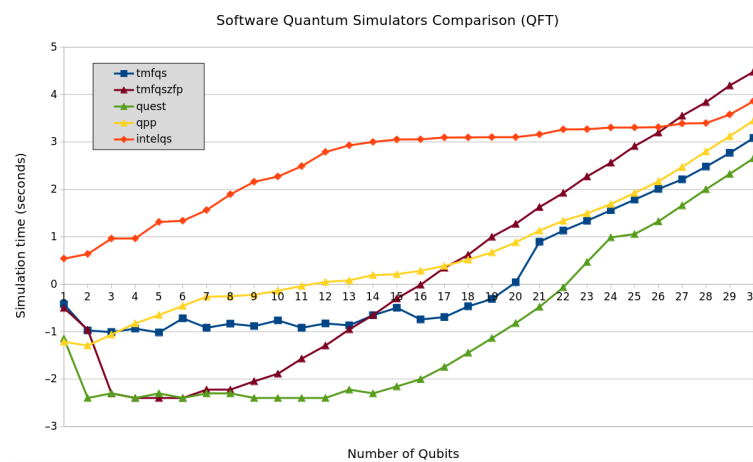


Figure 18. Quantum simulator comparison.

As can be seen in the graph in Figure 18, the Intel-QS simulator performs lower than the other simulators. QuEST exhibits the best performance. TMFQS performs acceptably

compared to these mature tools that have been optimized, for example, by using libraries such as MKL in the case of Intel-QS.

## 6. Conclusions

In conclusion, building a software quantum simulator requires a delicate balance between theoretical understanding and practical implementation strategies. The limitations of current quantum hardware, including qubit count and quality, drive the need for quantum simulators that allow researchers to explore quantum algorithms on classical computers. This work has shown that memory management techniques, such as dynamic pruning, full-state representation, and data compression, are essential for optimizing the simulation of quantum systems. While pruning techniques introduce certain challenges, such as fidelity loss and increased computational complexity, full-state representation with parallelization (via OpenMP or MPI) provides a robust framework for simulating larger quantum states. The use of data compression, such as ZFP, further extends the capacity to simulate a greater number of qubits without exceeding memory limits, though it introduces some overhead in processing time.

The comparative performance of the prototype simulator against established simulators like Intel-QS, QuEST, and qsim demonstrates the viability of these memory management techniques. By combining distributed and shared memory models, along with data compression, the simulator can handle increasingly complex simulations. Ultimately, this work contributes valuable insights into making quantum computing simulations more scalable and efficient, supporting the broader field of quantum computing as it advances toward practical applications.

### *Further Work*

While this work has explored several aspects of building software quantum computing simulators, there are numerous areas for further research.

- **Advanced Memory Compression Techniques:** Investigating more sophisticated data compression algorithms tailored specifically for quantum state vectors, and exploring hybrid compression techniques, combining lossy and lossless methods to balance between accuracy and memory usage.
- **Optimization of Quantum Gate Application:** Developing efficient algorithms for applying higher-order multi-qubit gates.
- **Error Mitigation and Fault Tolerance:** Developing strategies for incorporating error mitigation techniques within the simulators to emulate real quantum hardware better.
- **Distributed Quantum Computing:** Developing protocols for efficient communication and synchronization between distributed nodes to minimize overhead.
- **Energy-efficient Quantum Simulations:** Researching into energy-efficient memory management techniques to reduce the power consumption of large-scale quantum simulations.

By addressing these topics, future research can further enhance the efficiency, scalability, and the practicality of software quantum computing simulators, bringing them closer to the capabilities required for real-world quantum computing applications.

**Author Contributions:** G.D.: Project development, conceptualization, investigation, formal analysis, software development, methodology, simulations executions, manuscript writing. L.S.: supervision and review. C.B.: supervision and review. J.C.: supervision and review. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research received no external funding.

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** This work did not utilize any input data or generate new output data for analysis. Instead, simulations were conducted, and their results were compared with those from other leading simulators in the field. The outcomes are fully reproducible by running the example scenarios provided with these simulators.

**Acknowledgments:** This work used the high-performance computing resources of the Supercomputing and Scientific Computing Center (SC3) at the Universidad Industrial de Santander and the HPC facilities of the Université de Reims Champagne-Ardenne. We gratefully acknowledge the support and assistance of the system administrators and technical staff at both institutions. This study was presented at 11th Latin American High Performance Computing Conference, CARLA 2024, Santiago de Chile, Chile, 30 September–4 October 2024.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. Report, Q.C. Qbit Count. 2019. Available online: <https://quantumcomputingreport.com/scorecards/qubit-count/> (accessed on 12 June 2025).
2. Quantiki. List of QC Simulators. 2019. Available online: <https://www.quantiki.org/wiki/list-qc-simulators> (accessed on 12 June 2025).
3. Fingerhuth, M. Open-Source Quantum Software Projects. 2019. Available online: [https://github.com/qosf/os\\_quantum\\_software](https://github.com/qosf/os_quantum_software) (accessed on 12 June 2025).
4. Team, Q.O.S.F. Quantum Open Source Foundation. 2019. Available online: <https://qosf.org/> (accessed on 12 June 2025).
5. Nielsen, M.A.; Chuang, I.L. *Quantum Computation and Quantum Information*, 10th anniversary ed.; Cambridge University Press: Cambridge, UK, 2010.
6. Holevo, A.S. Bounds for the quantity of information transmitted by a quantum communication channel. *Probl. Inf. Transm.* **1973**, *9*, 177–183.
7. Williams, C.P. *Explorations in Quantum Computing*, 2nd ed.; Texts in Computer Science; Springer: Berlin/Heidelberg, Germany, 2011. [[CrossRef](#)]
8. Eleanor, R.; Wolfgang, P. *Quantum Computing, A Gentle Introduction*; The MIT Press: Cambridge, MA, USA, 2011.
9. Bergou, J.A.; Hillery, M. *Introduction to the Theory of Quantum Information Processing*; Springer Publishing Company, Incorporated: Berlin/Heidelberg, Germany, 2013.
10. Horodecki, R.; Horodecki, P.; Horodecki, M.; Horodecki, K. Quantum entanglement. *Rev. Mod. Phys.* **2009**, *81*, 865–942. [[CrossRef](#)]
11. Preskill, J. Lecture Notes for Physics 229: Quantum Information and Computation. 1998. Available online: [https://www.lorentz.leidenuniv.nl/quantumcomputers/literature/preskill\\_1\\_to\\_6.pdf](https://www.lorentz.leidenuniv.nl/quantumcomputers/literature/preskill_1_to_6.pdf) (accessed on 12 June 2025).
12. Häner, T.; Steiger, D.S.; Reiher, M.; Troyer, M. A low-depth quantum simulation of materials. *Npj Quantum Inf.* **2022**, *8*, 1–7. [[CrossRef](#)]
13. Villalonga, B.; Lyakh, D.; Boixo, S.; Neven, H.; Humble, T.S.; Biswas, R.; Rieffel, E.G.; Ho, A.; Mandrà, S. Establishing the quantum supremacy frontier with a 281 PFlop/s simulation. *Quantum Sci. Technol.* **2020**, *5*, 034003. [[CrossRef](#)]
14. Ding, Y.; Zhao, Z.; Li, A.; Chen, J.Z.; Li, D.; Zhang, H.; Tang, Y.; Chen, Y.; Chen, S. Sparse state simulation of quantum circuits. *ACM Trans. Quantum Comput.* **2023**, *4*, 1–25. [[CrossRef](#)]
15. Wu, J.; Li, A.; Ding, Y.; Li, D. High-performance quantum circuit simulation on GPUs. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '23), Denver, CO, USA, 12–17 November 2023; pp. 1–13. [[CrossRef](#)]
16. Insider, Q. Top 63 Quantum Computer Simulators for 2024. 2024. Available online: <https://thequantuminsider.com> (accessed on 12 May 2024).
17. Guerreschi, G.G.; Hogaboam, J.; Baruffa, F.; Sawaya, N. Intel Quantum Simulator: A cloud-ready high-performance simulator of quantum circuits. *arXiv* **2020**. [[CrossRef](#)]
18. Gheorghiu, V. Quantum++: A modern C++ quantum computing library. *arXiv* **2014**, arXiv:1412.4704. [[CrossRef](#)] [[PubMed](#)]
19. Qsim. Available online: (accessed on 12 June 2025). [[CrossRef](#)]
20. Jones, T.; Brown, A.; Bush, I.; Benjamin, S.C. QuEST and High Performance Simulation of Quantum Computers. *Sci. Rep.* **2019**, *9*, 10736. [[CrossRef](#)] [[PubMed](#)]
21. Strano, D. Qrack. 2019. Available online: <https://vm6502q.readthedocs.io/en/latest/> (accessed on 12 June 2025).
22. Strano, D.; Bollay, B.; Blaauw, A.; Shammah, N.; Zeng, W.J.; Mari, A. Exact and approximate simulation of large quantum circuits on a single GPU. *arXiv* **2023**. [[CrossRef](#)]

23. Díaz, G. Prototype Quantum Computing Simulator. 2024. Available online: <https://github.com/diaztoro/TMFQSfullstate.git> (accessed on 12 June 2025)
24. Trieu, D.B. Large-Scale Simulations of Error-Prone Quantum Computation Devices. Ph.D Thesis, Universität Wuppertal, Jülich, Germany, 2009.
25. Smelyanskiy, M.; Sawaya, N.P.D.; Aspuru-Guzik, A. qHiPSTER: The Quantum High Performance Software Testing Environment. *arXiv* **2016**, arXiv:1601.07195.
26. Lindstrom, P. Fixed-Rate Compressed Floating-Point Arrays. *IEEE Trans. Vis. Comput. Graph.* **2014**, *20*, 2674–2683. [[CrossRef](#)] [[PubMed](#)]

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.