

Parallelization of Quantum State Vector Simulator

Atit Pokharel*, Nabin K. Singh*,

*Department of Electrical and Computer Engineering, UAH

1 Introduction

Advancements in quantum computing have accelerated in recent years at an unprecedented rate. Certain classes of problems are unsolvable on classical computers, like quantum chemistry, multi-variable optimization, and specific machine-learning tasks [1, 2]. Progress in theory and algorithms is rapid because of a large research community and industrial involvement. Despite this huge leap forward in research papers and algorithms, practical realization of large-scale and fault-tolerant quantum computing platforms has not reached [3]. As a result, researchers rely on quantum state vector simulators. They are mathematical tools that allow the study of quantum algorithms and quantum circuit behavior using classical computing resources [4].

A quantum state vector simulator represents the quantum state of an n -qubit system as a complex vector of size 2^n . This representation is supported by the Copenhagen interpretation [5] and enables tracing of the exact state evolution in a quantum circuit. Numerous simulators have been developed over the years, including Qiskit Aer [6] from IBM, PennyLane [7], TorchQuantum from MIT [8], Cirq by Google [9], etc. However, these simulators have a fundamental scalability challenge. For example, simulating a mere 30 qubits requires storing and processing over a billion complex numbers. Both the memory usage and computational cost grow exponentially with the number of qubits and circuit depth. For example, simulating a mere 30 qubits requires storing and processing over a billion complex numbers. Even shallow and moderate circuit sizes' simulation takes a tremendous time on a classical computer. With single-threaded execution, iterating through these massive arrays to calculate matrix products becomes extremely slow. This latency limits the ability of researchers to develop complex and useful circuit designs and validate them.

This project explores the use of OpenMP [10] for parallel computing to accelerate quantum state-vector simulation on multi-core CPUs. As mentioned, quantum states of an n -qubit system can be represented as complex state vectors. In these vectors, the actions of quantum gates mean simply some linear transformation derived from Pauli matrices [11]. Thus, we first define a state vector simulator in C++ that is able to simulate single-qubit and two-qubit gates through direct matrix operations. Then, as most of these tasks are matrix operations, the simulator is parallelized using OpenMP by exploiting the independent and data-parallel structure of gate applications. The independent subsets of state amplitudes can be calculated and updated in parallel.

1.1 Goals and Objectives

The main objective of this project is to bridge the gap between computationally intensive quantum simulations and efficient parallel computing resources. They are summarized as follows

- **Runtime Reduction:** Accelerate the computation of quantum state vector simulations by parallelizing the resource-intensive gate operations using OpenMP.
- **Modular Framework Development:** Design a modular simulation tool that easily allows the addition of new quantum gates and different circuit structures.
- **Scalability Analysis:** Evaluate how the implementation scales with an increasing number of qubits by profiling Execution time (T_p), Speedup (S_p), Cost (C_p), and Parallel Efficiency (E_p).
- **Performance Benchmarking:** Demonstrate the performance gains by comparing the above metrics of the parallelized approach against serial methods on multi-core architecture.

2 Primers on Quantum Computing and its Simulation

To understand the computational load of a simulator, we define the mathematical representation of quantum states and operations in this section.

Qubits and State Vectors: Anyone in computer-system academia or industry must have heard multiple times that "a classical bit exists in a state of 0 or 1, whereas a qubit exists in a superposition of both." If we write this mathematically, the state of a single qubit is represented by a column vector of two complex numbers, $|\psi\rangle = \alpha|0\rangle + \beta|1\rangle$, where $|\alpha|^2 + |\beta|^2 = 1$. For a system of N qubits, such a state space grows exponentially. The combined state is represented by the tensor product of individual qubit states, resulting in a vector of size 2^N , as

$$|\Psi\rangle = \alpha_0|00\dots 0\rangle + \alpha_1|00\dots 1\rangle + \dots + \alpha_{2^N-1}|11\dots 1\rangle$$

Pauli Matrices and Fundamental Gates: The building blocks of quantum operations are the Pauli matrices. These are 2×2 complex matrices that form a basis for the space of single-qubit operations. In this project, we explicitly implement the following gates, as

- **Identity (I):** Typical identity matrix, leaves the qubit unchanged.

$$I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

- **Pauli-X (X):** Similar to the classical NOT gate (also called bit flip).

$$X = \sigma_x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

- **Pauli-Y (Y):** Combines bit flip and phase flip.

$$Y = \sigma_y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}$$

- **Pauli-Z (Z):** Flips the phase of the $|1\rangle$ component.

$$Z = \sigma_z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

- **Hadamard (H):** Creates superposition from basis states.

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

In general, we represent them as U , since they are unitary. These gates acts on a state vector $|\psi\rangle$ and the new state $|\psi'\rangle$ is calculated via matrix-vector multiplication:

$$|\psi'\rangle = U|\psi\rangle$$

An example: The main challenge in simulation arises when applying operations to multi-qubit systems. If a gate acts on specific qubits while keeping others unchanged (like the Identity matrix I), the process is represented using the tensor product (\otimes)(Kronecker product). For example, applying a Pauli-X gate to the first qubit of a two-qubit system can be represented as:

$$U_{total} = X \otimes I = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \otimes \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \quad (1)$$

Computing these products allows for parallelization because the individual elements of the resulting large matrix can be computed independently.

2.1 A benchmark for us: Variational Quantum Circuit

To evaluate the performance of our parallel simulator, we choose Variational Quantum Circuits (VQCs) as a benchmark. VQCs are a type of quantum algorithm used widely in Quantum Machine Learning (QML) [12] and optimization tasks (e.g., VQE, QAOA). They are also referred to as Parameterized Quantum Circuits (PQCs).

A typical VQC consists of a fixed gate structure (the ansatz). Here, it is essential to notice that the above-mentioned gate matrices can also be represented as the complex exponentiation of Pauli matrices ($R_P(\theta) = e^{-i\theta P/2}$).

$$R_x(\theta) = \begin{pmatrix} \cos \frac{\theta}{2} & -i \sin \frac{\theta}{2} \\ -i \sin \frac{\theta}{2} & \cos \frac{\theta}{2} \end{pmatrix}$$

$$R_y(\theta) = \begin{pmatrix} \cos \frac{\theta}{2} & -\sin \frac{\theta}{2} \\ \sin \frac{\theta}{2} & \cos \frac{\theta}{2} \end{pmatrix}$$

$$R_z(\theta) = \begin{pmatrix} e^{-i\theta/2} & 0 \\ 0 & e^{i\theta/2} \end{pmatrix}$$

Here, the gate parameters (i.e., the rotation angles θ) are tunable/trainable parameters and the goal of a VQC is to optimize these θ values to approximate some function (similar to a neural network in the classical domain). Such a circuit typically alternates between:

1. **Rotation Layers:** Single-qubit rotation gates (e.g., $R_x(\theta)$, $R_y(\theta)$) applied to all qubits.
2. **Entanglement Layers (CNOT):** A fundamental two-qubit gate used to create entanglement. In Fig. 1, the two-qubit connection with an \oplus connection is the CNOT gate operation. It flips the target qubit if and only if the control qubit is $|1\rangle$.

$$\text{CNOT} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

Fig. 1 shows a typical structure of a VQC, where a fully entangled 5-qubit VQC with three arbitrary layers is shown.

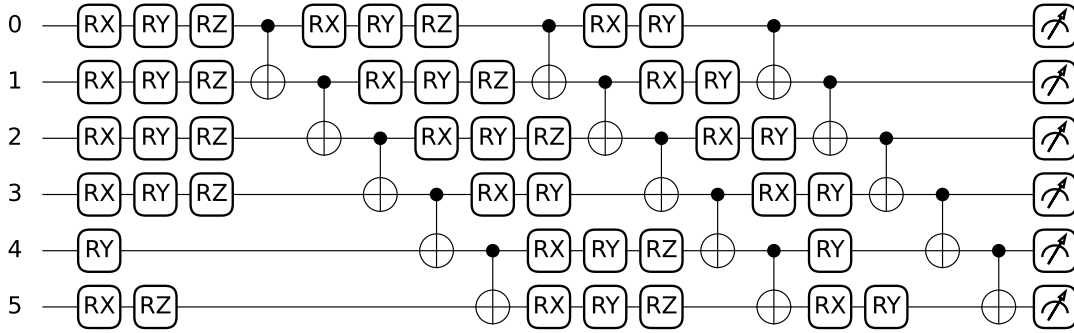


Figure 1: A typical 5-qubit PQC quantum circuit with three layers.

Why VQCs are an ideal benchmark: VQCs provide a solid test for simulation performance because they require the repeated application of dense matrix product operations across all qubits. Unlike simple algorithms that might only use a few gates, VQCs possess a deep and repetitive layered structure ("depth"). This stresses the simulator to compute complex state evolutions continuously. Such a structure allows us to easily scale the workload by simply increasing the number of qubits or the number of layers in the circuit. Thus, stress-testing the parallel efficiency of our OpenMP implementation becomes feasible and simple.

3 Methods

3.1 Theoretical Basis

The computational core of our simulation relies on linear algebra, particularly the Kronecker (tensor) product. Given an operator A of size $m \times n$ and an operator B of size $p \times q$, the product $C = A \otimes B$ results in a matrix of size $mp \times nq$. The element-wise expression is

$$C_{i,j} = A_{[i/p],[j/q]} \times B_{i \pmod{p}, j \pmod{q}}$$

Conceptually, a quantum circuit is processed layer by layer. If we were to construct the full operator exactly, for a gate U applied to qubit k , the global unitary would be:

$$U_{global} = I^{\otimes k} \otimes U \otimes I^{\otimes (N-k-1)}$$

However, naively constructing this U_{global} (a $2^N \times 2^N$ matrix) is memory-intensive (usually results in segmentation faults). Therefore, we implement a State-Vector Update approach. It applies the mathematical equivalent of these tensor products directly to the state amplitudes without forming those large, sparse matrices.

3.2 Implementation Environment

The simulator was implemented in C++ (C++17) as it handles low-level memory management. Parallelization was achieved using the OpenMP API, and the complex number arithmetic was handled using the standard `std::complex<double>` library.

3.3 Simulation Flow

The simulation proceeds in the following steps:

1. **Initialization:** The state vector $|\Psi\rangle$ is initialized to the ground state $|0 \dots 0\rangle$ (a vector of size 2^N with the first element as 1 and all others 0).
2. **Layer Execution:** The circuit is processed sequentially. For each gate, we identify the target qubits and update the state vector in place.
3. **State Evolution:** The operations transform the complex amplitudes of the state vector (linear algebra based), effectively evolving the state of the system in Hilbert space.

3.4 Parallel Gate Application Algorithm

For a single-qubit gate acting on a target qubit t , we iterate through the state vector to identify pairs of indices (i_0, i_1) where the t -th bit is 0 and 1, respectively. We then update these amplitudes using the 2×2 gate matrix.

OpenMP Parallelization Strategy: We used the `#pragma omp parallel for collapse(2)` directive. This instructs the compiler to:

1. Spawn a team of threads.
2. Collapse the nested loops (iterating over the state vector blocks) into a single iteration space.
3. Distribute these iterations dynamically across the available CPU cores.

Listing 1: Parallel application of a single-qubit gate using OpenMP

```
void apply_single_qubit_gate(State &state, int n_qubits,
                           int target, const Mat2 &U) {
    const std::size_t dim = state.size();
    const std::size_t step = (1ULL << target);
    const std::size_t period = step * 2;

    // PARALLEL REGION: Automatically distributes loop iterations
    #pragma omp parallel for collapse(2)
    for (std::size_t base = 0; base < dim; base += period) {
        for (std::size_t i = 0; i < step; ++i) {
            std::size_t i0 = base + i;
            std::size_t i1 = base + i + step;

            Complex a0 = state[i0];
            Complex a1 = state[i1];

            // In-place update using the 2x2 Unitary U
            state[i0] = U.m[0][0] * a0 + U.m[0][1] * a1;
            state[i1] = U.m[1][0] * a0 + U.m[1][1] * a1;
        }
    }
}
```

3.5 Entangling Gates (CNOT)

The CNOT gate flips the target qubit if and only if the control qubit is $|1\rangle$. In our state-vector representation, this is done by swapping the amplitudes at indices where the control bit is 1 and the target bit differs. This operation is memory-bound rather than compute-bound. OpenMP is used here to parallelize the traversal of the state array.

Listing 2: Parallel implementation of the CNOT gate

```
void apply_cx(State &state, int n_qubits, int control, int target) {
    const std::size_t dim = state.size();
    std::size_t control_mask = (1ULL << control);
    std::size_t target_mask = (1ULL << target);

    // PARALLEL REGION
    #pragma omp parallel for
    for (std::size_t idx = 0; idx < dim; ++idx) {
```

```

// Only act if control bit is 1
if ( (idx & control_mask) != 0 ) {
    std::size_t flipped_idx = idx ^ target_mask;

    // Swap only once per pair to avoid race conditions
    if (idx < flipped_idx) {
        std::swap(state[idx], state[flipped_idx]);
    }
}
}
}

```

4 Experiments and Results

For the evaluation of our parallel quantum circuit simulator, a series of experiments were conducted with the Variational Quantum Circuit (VQC) benchmark, described in Section II. The experiments were conducted with objective to measure the impact of OpenMP parallelization on the parameters: execution time (T_p), Speedup (S_p), Parallel Efficiency (E_p), and Computational Cost (C_p).

Simulation Settings: All experiments were conducted on a **Lambda Vector One** high-performance deep learning workstation, with the following specific hardware specifications:

- **Processor:** AMD Ryzen™ 9 7950X (Zen 4 architecture).
- **Core Count:** 16 physical cores, 32 logical threads (Simultaneous Multithreading).
- **Memory:** 64 GB DDR5 RAM.
- **Operating System:** Ubuntu Linux (via Lambda Stack).

This hardware configuration enabled us to test explicit scaling from 1 thread to the maximum available 32 logical threads.

4.1 Test Case 1: Scaling with Circuit Depth

For the first experiment, the system size was kept constant(20 qubits), and tested for the different number of layers in the VQC from 3 to 8. Thus, this tests the simulator’s ability to handle deeper circuits where the state vector is continuously updated.

Figure 2 illustrates the execution time as a function of thread count. In the experiment we observed the significant decrease in runtime on the increase of number of threads. For a depth of 8 layers, the serial execution time was 12.11 seconds, but with 15 threads, this dropped to 1.72 seconds, showcasing a substantial performance gain. However, beyond 15 threads, the stagnation in performance was observed. This indicates that for 20 qubits, the workload per thread becomes low compared to the thread management overhead and memory bandwidth limitation, which restrict the further speedup.

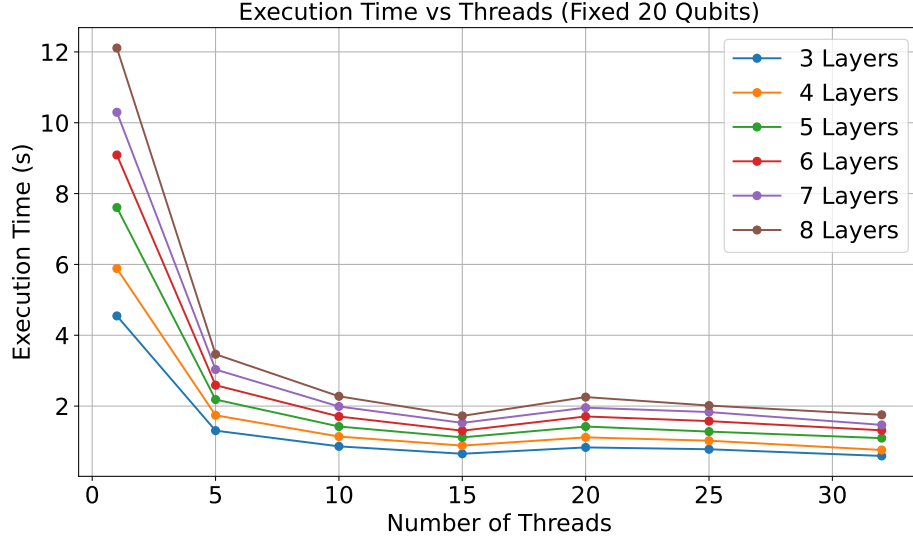


Figure 2: Execution Time vs. Number of Threads for fixed 20 qubits with varying circuit depths (3 to 8 layers).

4.2 Test Case 2: Scaling with System Size

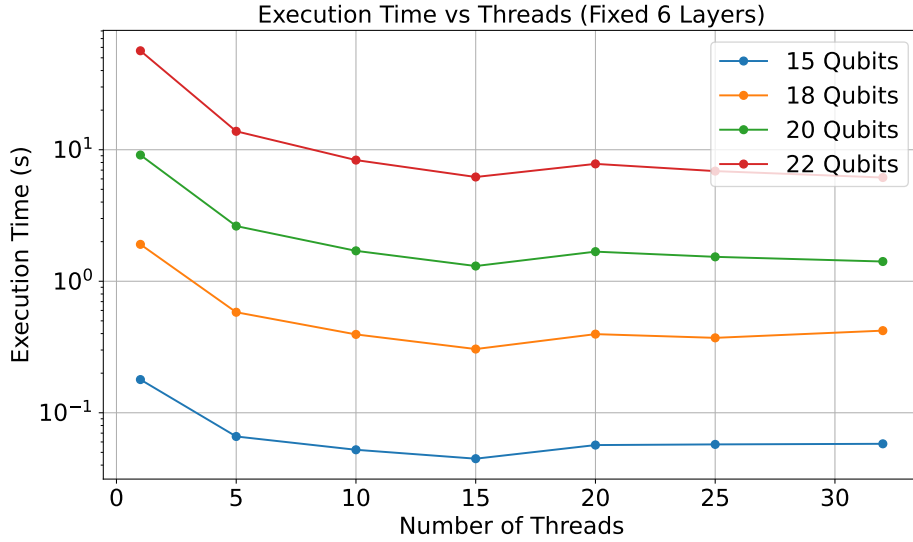


Figure 3: Execution Time vs. Threads for varying qubit counts (Log Scale). Note the exponential increase in serial time as qubits increase.

for the second experiment, we made the circuit depth constant (to 6 layers) and tested for the different number of qubits (from 5 to 22). This is the most critical test, as the state vector size grows exponentially (2^N).

Small System Overhead: For smaller systems (5 and 10 qubits), negligible or negative results (see data for 5 qubits) were observed with parallelization. The overhead of spawning

threads and managing parallel regions outweighed the computational cost of the trivial state vector updates.

Large System Acceleration: However for larger systems, the benefits of parallelization were evident. Figure 3 shows the execution time on a logarithmic scale.

- **22 Qubits:** The serial execution took **56.55 seconds**. Using 15 threads reduced this to **6.20 seconds**, a $\approx 9\times$ improvement.
- **20 Qubits:** Execution time dropped from 9.12 seconds (serial) to 1.30 seconds (15 threads).

4.3 Performance Metrics Analysis

Speedup (S_p):

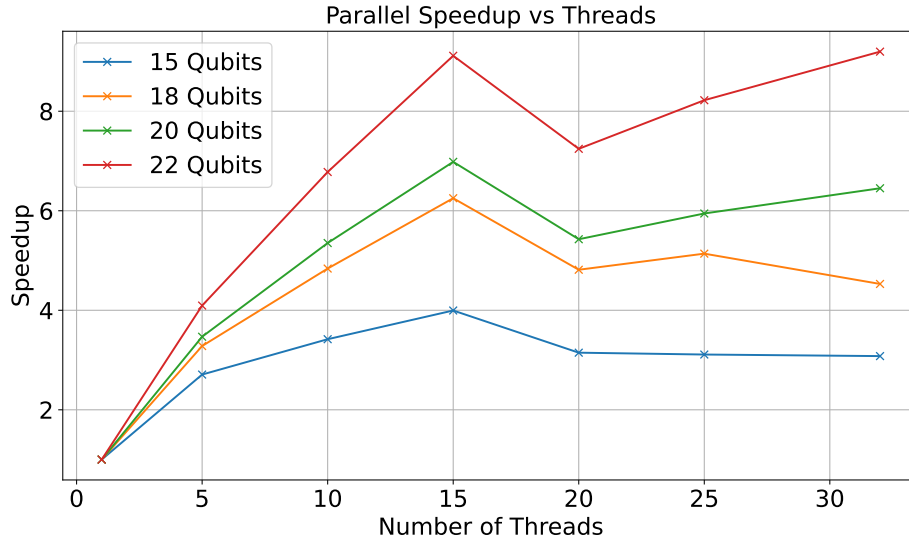


Figure 4: Speedup (S_p) vs. Threads. Larger qubit systems exhibit better speedup characteristics.

Speedup is defined as $S_p = T_1/T_p$. Figure 4 shows the speedup achieved for varying Qubit counts. The plot clearly demonstrates that larger workloads (22 Qubits) achieve better speedup because of the higher ratio of useful computation to parallel overhead. We observed a maximum speedup of approximately $9.1\times$ at 22 Qubits.

Parallel Efficiency (E_p):

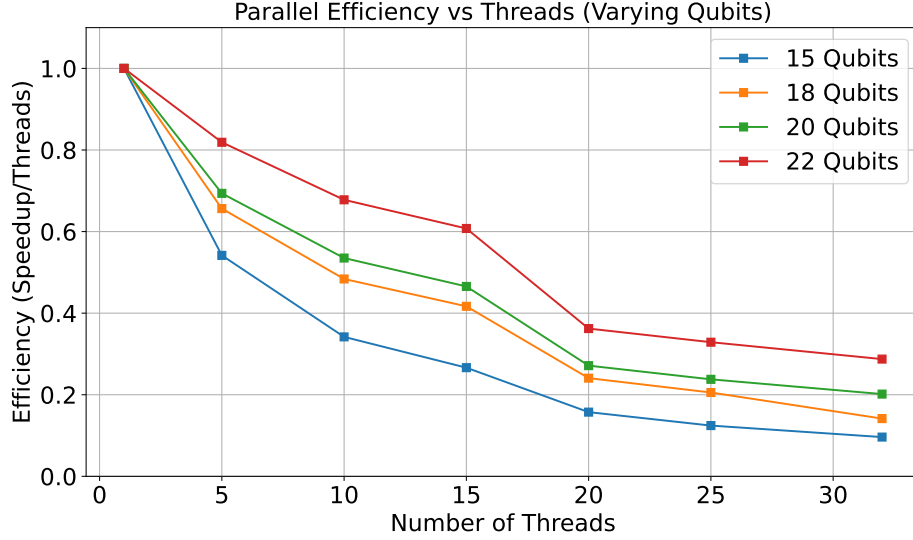


Figure 5: Parallel Efficiency vs. Threads. The drop in efficiency at higher thread counts indicates memory bandwidth saturation.

Efficiency ($E_p = S_p/P$), measures how effectively the utilized cores are being used. Figure 5 demonstrates, how efficiency drops as threads increases. The state-vector update method is highly memory-intensive (traversing large arrays). As we increase the number of threads, we eventually saturate the memory bandwidth of the hardware. This explains why the efficiency drops below 50% when the thread count exceeds 15, and why adding threads beyond this point (e.g., 32 threads) doesn't yield any significant results.

Computational Cost:

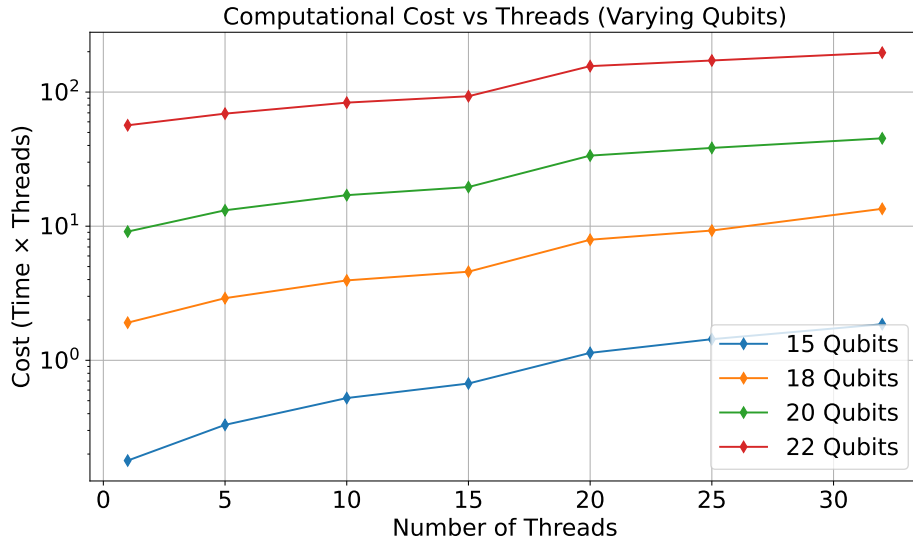


Figure 6: Computational Cost (Time \times Threads) vs. Threads. The upward trend represents the overhead introduced by parallelization.

We also analyzed the Computational Cost ($\text{Time} \times \text{Threads}$). Ideally, cost should remain constant. However, as shown in Figure 6, the cost increases with the number of threads. This increase reflects the “price” of parallelism, where the total CPU cycles grow due to parallel overhead, even though the wall-clock execution time decreases.

5 Conclusion

This project demonstrates the use of parallelization (OpenMP) as an effective way to accelerate quantum state vector simulation on a classical multi-core CPU. We simulated quantum computing in parallel on a classical computer because quantum hardware is not yet easily accessible. Furthermore, executing such simulations on a classical computer using a single core takes an extremely long time, since the state vector grows exponentially as the number of Qubits increases. Therefore, by running single-qubit and two-qubit gate operations in a C++ simulator across multiple CPU cores, the workload was distributed among the cores, which significantly reduced the execution time.

The experiments demonstrated an increase in speedup for larger quantum systems as the number of threads increased. However, the performance gains eventually stagnated after a certain threshold of threads due to parallel overhead and memory bandwidth limitations, which prevented further speedup for the same number of Qubits (state vectors). Overall, this work confirms that OpenMP-based parallel computing improves quantum simulations mainly for larger workloads and up to a certain number of threads, while overhead and limitation in memory bandwidth, limit the further gain in the performance.

References

- [1] R. Courtland, “Google aims for quantum computing supremacy [news],” *IEEE Spectrum*, vol. 54, no. 6, pp. 9–10, 2017.
- [2] C. P. Koch, U. Boscain, T. Calarco, G. Dirr, S. Filipp, S. J. Glaser, R. Kosloff, S. Montangero, T. Schulte-Herbrüggen, D. Sugny *et al.*, “Quantum optimal control in quantum technologies. strategic report on current status, visions and goals for research in europe,” *EPJ Quantum Technology*, vol. 9, no. 1, p. 19, 2022.
- [3] S. Resch and U. R. Karpuzcu, “Benchmarking quantum computers and the impact of quantum noise,” *ACM Computing Surveys (CSUR)*, vol. 54, no. 7, pp. 1–35, 2021.
- [4] S. Endo, J. Sun, Y. Li, S. C. Benjamin, and X. Yuan, “Variational quantum simulation of general processes,” *Physical Review Letters*, vol. 125, no. 1, p. 010501, 2020.
- [5] H. P. Stapp, “The copenhagen interpretation,” *American journal of physics*, vol. 40, no. 8, pp. 1098–1116, 1972.
- [6] J. Faj, I. Peng, J. Wahlgren, and S. Markidis, “Quantum computer simulations at warp speed: Assessing the impact of gpu acceleration: A case study with ibm qiskit aer, nvidia thrust & cuquantum,” in *2023 IEEE 19th International Conference on e-Science (e-Science)*. IEEE, 2023, pp. 1–10.

- [7] V. Bergholm, J. Izaac, M. Schuld, C. Gogolin, S. Ahmed, V. Ajith, M. S. Alam, G. Alonso-Linaje, B. AkashNarayanan, A. Asadi *et al.*, “PennyLane: Automatic differentiation of hybrid quantum-classical computations,” *arXiv preprint arXiv:1811.04968*, 2018.
- [8] H. Wang, Z. Liang, J. Gu, Z. Li, Y. Ding, W. Jiang, Y. Shi, D. Z. Pan, F. T. Chong, and S. Han, “Torchquantum case study for robust quantum circuits,” in *Proceedings of the 41st IEEE/ACM International Conference on Computer-Aided Design*, 2022, pp. 1–9.
- [9] S. V. Isakov, D. Kafri, O. Martin, C. V. Heidweiller, W. Mroczkiewicz, M. P. Harrigan, N. C. Rubin, R. Thomson, M. Broughton, K. Kissell *et al.*, “Simulations of quantum circuits with approximate noise using qsim and cirq,” *arXiv preprint arXiv:2111.02396*, 2021.
- [10] L. Dagum and R. Menon, “Openmp: an industry standard api for shared-memory programming,” *IEEE computational science and engineering*, vol. 5, no. 1, pp. 46–55, 1998.
- [11] D. Gottesman, “The heisenberg representation of quantum computers,” *arXiv preprint quant-ph/9807006*, 1998.
- [12] A. Pokharel, R. Rahman, T. Morris, and D. C. Nguyen, “Quantum federated learning for multimodal data: A modality-agnostic approach,” in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 545–554.

A Appendix

The source code is presented in this section. There are two C++ files: one defines the quantum circuit and associated utilities, and the other runs the benchmark. The .h header file is also provided.

The latter part consists of the bash script to run the whole project and its output file.

quantum_circuit.cpp

```

1  #include "quantum_circuit.h"
2  #include <iostream>
3  #include <cmath>
4  #include <iomanip>
5  #include <omp.h>
6
7  //Gate Definitions
8  Mat2 make_I() {
9      return Mat2{ { { Complex(1,0), Complex(0,0) },
10                     { Complex(0,0), Complex(1,0) } } };
11 }
12
13 Mat2 make_X() {

```

```

14     return Mat2{ { { Complex(0,0), Complex(1,0) },
15                   { Complex(1,0), Complex(0,0) } } };
16 }
17
18 Mat2 make_Y() {
19     return Mat2{ { { Complex(0,0), Complex(0,-1) },
20                   { Complex(0,1), Complex(0, 0) } } };
21 }
22
23 Mat2 make_Z() {
24     return Mat2{ { { Complex(1,0), Complex(0,0) },
25                   { Complex(0,0), Complex(-1,0) } } };
26 }
27
28 Mat2 make_H() {
29     const double inv_sqrt2 = 1.0 / std::sqrt(2.0);
30     return Mat2{ { { Complex(inv_sqrt2,0), Complex(inv_sqrt2,0) },
31                   { Complex(inv_sqrt2,0), Complex(-inv_sqrt2,0) } } };
32 }
33 //Rotational Gates
34 Mat2 make_RX(double theta) {
35     double c = std::cos(theta / 2.0);
36     double s = std::sin(theta / 2.0);
37     return Mat2{ { { Complex(c, 0), Complex(0, -s) },
38                   { Complex(0, -s), Complex(c, 0) } } };
39 }
40
41 Mat2 make_RY(double theta) {
42     double c = std::cos(theta / 2.0);
43     double s = std::sin(theta / 2.0);
44     return Mat2{ { { Complex(c, 0), Complex(-s, 0) },
45                   { Complex(s, 0), Complex(c, 0) } } };
46 }
47
48 Mat2 make_RZ(double theta) {
49     double c = std::cos(theta / 2.0);
50     double s = std::sin(theta / 2.0);
51     //  $e^{-i\theta/2} = \cos - i \sin$ 
52     //  $e^{+i\theta/2} = \cos + i \sin$ 
53     return Mat2{ { { Complex(c, -s), Complex(0, 0) },
54                   { Complex(0, 0), Complex(c, s) } } };
55 }
56
57
58 //Operation Functions
59 State make_zero_state(int n_qubits) {
60     std::size_t dim = (1ULL << n_qubits);
61     State psi(dim);
62

```

```

63     // Parallel initialization for large arrays
64     #pragma omp parallel for
65     for(size_t i=0; i<dim; ++i) {
66         psi[i] = Complex(0,0);
67     }
68
69     psi[0] = Complex(1.0, 0.0);
70     return psi;
71 }
72
73 void apply_single_qubit_gate(State &state, int n_qubits, int target, const Mat2 &U)
74 ↪ {
75     const std::size_t dim = state.size();
76     // validation check
77     if (dim != (1ULL << n_qubits)) return;
78
79     const std::size_t step = (1ULL << target);
80     const std::size_t period = step * 2;
81
82     // PARALLEL REGION:
83     #pragma omp parallel for collapse(2)
84     for (std::size_t base = 0; base < dim; base += period) {
85         for (std::size_t i = 0; i < step; ++i) {
86             std::size_t i0 = base + i;
87             std::size_t i1 = base + i + step;
88
89             Complex a0 = state[i0];
90             Complex a1 = state[i1];
91
92             // Matrix-Vector multiplication for single qubit gate
93             Complex b0 = U.m[0][0] * a0 + U.m[0][1] * a1;
94             Complex b1 = U.m[1][0] * a0 + U.m[1][1] * a1;
95
96             state[i0] = b0;
97             state[i1] = b1;
98         }
99     }
100
101 void apply_cx(State &state, int n_qubits, int control, int target) {
102     const std::size_t dim = state.size();
103
104     if (control == target) return;
105
106     std::size_t control_mask = (1ULL << control);
107     std::size_t target_mask = (1ULL << target);
108
109     // PARALLEL REGION:
110     #pragma omp parallel for

```

```

111     for (std::size_t idx = 0; idx < dim; ++idx) {
112         if ( (idx & control_mask) != 0 ) {
113             std::size_t flipped_idx = idx ^ target_mask;
114             // Ensure we swap only once per pair
115             if (idx < flipped_idx) {
116                 std::swap(state[idx], state[flipped_idx]);
117             }
118         }
119     }
120 }
121
122 void print_probabilities(const State &state) {
123     const std::size_t dim = state.size();
124     if (dim > 32) {
125         std::cout << " (State vector too large to print: size=" << dim << ")\n";
126         return;
127     }
128
129     std::cout << std::fixed << std::setprecision(4);
130     for (std::size_t i = 0; i < dim; ++i) {
131         double p = std::norm(state[i]);
132         if (p > 1e-10) {
133             std::cout << "|" << i << "> : " << p << "\n";
134         }
135     }
136 }

```

quantum_circuit.h

```

1     #ifndef QUANTUM_CIRCUIT_H
2     #define QUANTUM_CIRCUIT_H
3
4     #include <complex>
5     #include <vector>
6
7     //type def
8     using Complex = std::complex<double>;
9     using State   = std::vector<Complex>;
10
11     struct Mat2 {
12         Complex m[2][2];
13     };
14
15     //Gate Functions
16     Mat2 make_I();
17     Mat2 make_X();
18     Mat2 make_Y();
19     Mat2 make_Z();
20     Mat2 make_H();

```

```

21 Mat2 make_RX(double theta);
22 Mat2 make_RY(double theta);
23 Mat2 make_RZ(double theta);
24
25 //Operation Functions
26 // Initialize state  $|0...0\rangle$  (zero state)
27 State make_zero_state(int n_qubits);
28
29 //single-qubit gate U to target qubit (single qubit gate)
30 void apply_single_qubit_gate(State &state, int n_qubits, int target, const Mat2 &U);
31
32 // CX (CNOT)
33 void apply_cx(State &state, int n_qubits, int control, int target);
34
35 //print probabilities
36 void print_probabilities(const State &state);
37
38 #endif
39
40

```

simple_quantum_circuit.cpp

```

1      #include "quantum_circuit.h"
2  #include <iostream>
3  #include <omp.h>
4  #include <vector>
5  #include <random>
6
7  //VQC Benchmark
8  void run_benchmark(int n_qubits, int layers, int training_rounds) {
9      std::cout << "\nVQC Benchmark (" << n_qubits << " qubits)\n";
10     std::cout << "Threads: " << omp_get_max_threads() << "\n";
11     std::cout << "Simulating " << training_rounds << " training rounds (forward
    ↳ passes only.)\n";
12
13     double total_start_time = omp_get_wtime();
14     std::mt19937 gen(42); //random number generator
15     std::uniform_real_distribution<> dis(0.0, 2.0 * M_PI); //uniform distribution for
    ↳ VQC parameters
16
17     for (int round = 0; round < training_rounds; ++round) {
18         double round_start = omp_get_wtime();
19         //reset state for each forward pass
20         State psi = make_zero_state(n_qubits);
21         for (int layer = 0; layer < layers; ++layer) { //for each layer
22             // Parametrized gate to all: RX
23             for(int i = 0; i < n_qubits; ++i) {
24                 double theta = dis(gen);

```



```

25         apply_single_qubit_gate(psi, n_qubits, i, make_RX(theta));
26     }
27     //Parametrized gate RY
28     for(int i = 0; i < n_qubits; ++i) {
29         double theta = dis(gen);
30         apply_single_qubit_gate(psi, n_qubits, i, make_RY(theta));
31     }
32     //RZ
33     for(int i = 0; i < n_qubits; ++i) {
34         double theta = dis(gen);
35         apply_single_qubit_gate(psi, n_qubits, i, make_RZ(theta));
36     }
37     //All to all entanglement
38     for(int i = 0; i < n_qubits; ++i) {
39         for(int j = 0; j < n_qubits; ++j) {
40             if (i != j) {
41                 apply_cx(psi, n_qubits, i, j);
42             }
43         }
44     }
45 }
46 double round_end = omp_get_wtime();
47 std::cout << "Round " << (round + 1) << "/" << training_rounds
48     << " complete. Time: " << (round_end - round_start) << " s\n";
49 }
50
51 double total_end_time = omp_get_wtime();
52 std::cout << "\nRun Complete.\n";
53 std::cout << "Total Time: " << (total_end_time - total_start_time) << "
54     ↪ seconds.\n";
55 std::cout << "Average Time per Round: " << (total_end_time - total_start_time) /
56     ↪ training_rounds << " seconds.\n";
57 }
58 //main
59 int main(int argc, char* argv[]) {
60     //default cases
61     int n_qubits = 20;
62     int layers = 6;
63     int rounds = 10;
64
65     if (argc > 1) n_qubits = std::atoi(argv[1]);
66     if (argc > 2) layers = std::atoi(argv[2]);
67     if (argc > 3) rounds = std::atoi(argv[3]);
68
69     run_benchmark(n_qubits, layers, rounds);
70
71     return 0;

```

```
72 }  
73
```

Bash Script:

```
1      #!/bin/bash  
2  #compile first  
3  echo "Compiling..."  
4  g++ -O3 -fopenmp simple_quantum_circuit.cpp quantum_circuit.cpp -o  
    ↪ simple_quantum_circuit  
5  OUTPUT_FILE="benchmark_results.txt" #output file  
6  #clear ir first if file exists  
7  echo "Results:" > $OUTPUT_FILE  
8  
9  THREADS_LIST=(1 5 10 15 20 25 32)  
10 ROUNDS=5 # Reduced for speed in automated runs, increase if needed  
11  
12 #test case 1  
13 QUBITS_FIXED=20  
14 LAYERS_LIST=(3 4 5 6 7 8)  
15  
16 echo "Test Case 1: Varying Layers (Fixed 20 Qubits)"  
17 echo "" >> $OUTPUT_FILE  
18 echo "TEST CASE 1: Fixed Qubits = $QUBITS_FIXED, Varying Layers" >> $OUTPUT_FILE  
19 printf "%-10s %-10s %-10s %-20s\n" "Threads" "Qubits" "Layers" "Total Time (s)" >>  
    ↪ $OUTPUT_FILE #header  
20  
21 for LAYERS in "${LAYERS_LIST[@]"; do #for each layer  
22     for THREADS in "${THREADS_LIST[@]"; do  
23         echo "Running: Qubits=$QUBITS_FIXED, Layers=$LAYERS, Threads=$THREADS"  
24         ↪ #running  
25  
26         RESULT=$(OMP_NUM_THREADS=$THREADS ./simple_quantum_circuit $QUBITS_FIXED  
27         ↪ $LAYERS $ROUNDS) #run benchmark  
28         #extract total time  
29         TIME=$(echo "$RESULT" | grep "Total Time:" | awk '{print $3}')  
30  
31         printf "%-10s %-10s %-10s %-20s\n" "$THREADS" "$QUBITS_FIXED" "$LAYERS"  
32         ↪ "$TIME" >> $OUTPUT_FILE  
33     done  
34     echo "" >> $OUTPUT_FILE # Spacer between layer groups  
35 done  
36  
37 #Test Case 2: Varying Qubits (Fixed 6 Layers)  
38 LAYERS_FIXED=6  
39 QUBITS_LIST=(5 10 15 18 20 22)  
40  
41 echo "Test Case 2: Varying Qubits (Fixed 6 Layers)"  
42 echo "" >> $OUTPUT_FILE
```

```

40 echo "TEST CASE 2: Fixed Layers = $LAYERS_FIXED, Varying Qubits" >> $OUTPUT_FILE
41 printf "%-10s %-10s %-10s %-20s\n" "Threads" "Qubits" "Layers" "Total Time (s)" >>
   ↳ $OUTPUT_FILE #header
42
43 for QUBITS in "${QUBITS_LIST[@]}; do #for each qubit
44     for THREADS in "${THREADS_LIST[@]}; do
45         echo "Running: Qubits=$QUBITS, Layers=$LAYERS_FIXED, Threads=$THREADS"
46
47         #running
48         RESULT=$(OMP_NUM_THREADS=$THREADS ./simple_quantum_circuit $QUBITS
   ↳ $LAYERS_FIXED $ROUNDS)
49
50         #extract total time
51         TIME=$(echo "$RESULT" | grep "Total Time:" | awk '{print $3}')
52
53         printf "%-10s %-10s %-10s %-20s\n" "$THREADS" "$QUBITS" "$LAYERS_FIXED"
   ↳ "$TIME" >> $OUTPUT_FILE
54     done
55     echo "" >> $OUTPUT_FILE
56 done
57
58 echo "Done. $OUTPUT_FILE"

```

Output (Run 1:)

```

1           Results:
2
3 TEST CASE 1: Fixed Qubits = 20, Varying Layers
4 Threads    Qubits    Layers    Total Time (s)
5 1           20         3         4.54596
6 5           20         3         1.30675
7 10          20         3         0.863134
8 15          20         3         0.653564
9 20          20         3         0.833679
10 25         20         3         0.781296
11 32         20         3         0.592721
12
13 1           20         4         5.88434
14 5           20         4         1.74037
15 10          20         4         1.14022
16 15          20         4         0.884248
17 20          20         4         1.11676
18 25          20         4         1.02337
19 32          20         4         0.760781
20
21 1           20         5         7.60661
22 5           20         5         2.18475
23 10          20         5         1.42377
24 15          20         5         1.11657

```

25	20	20	5	1.4235
26	25	20	5	1.28007
27	32	20	5	1.09518
28				
29	1	20	6	9.08939
30	5	20	6	2.58956
31	10	20	6	1.7066
32	15	20	6	1.30513
33	20	20	6	1.70739
34	25	20	6	1.57605
35	32	20	6	1.31584
36				
37	1	20	7	10.2956
38	5	20	7	3.03262
39	10	20	7	1.98872
40	15	20	7	1.52616
41	20	20	7	1.95523
42	25	20	7	1.8321
43	32	20	7	1.46896
44				
45	1	20	8	12.1098
46	5	20	8	3.46344
47	10	20	8	2.27673
48	15	20	8	1.72295
49	20	20	8	2.25464
50	25	20	8	2.01544
51	32	20	8	1.75398
52				
53				
54	TEST CASE 2: Fixed Layers = 6, Varying Qubits			
55	Threads	Qubits	Layers	Total Time (s)
56	1	5	6	0.000264091
57	5	5	6	0.00153651
58	10	5	6	0.00201087
59	15	5	6	0.00245149
60	20	5	6	0.00328084
61	25	5	6	0.00374911
62	32	5	6	0.00427548
63				
64	1	10	6	0.00388986
65	5	10	6	0.00698495
66	10	10	6	0.00770696
67	15	10	6	0.00843069
68	20	10	6	0.010934
69	25	10	6	0.0122837
70	32	10	6	0.014421
71				
72	1	15	6	0.178828
73	5	15	6	0.0660308

74	10	15	6	0.052311
75	15	15	6	0.0447584
76	20	15	6	0.0568088
77	25	15	6	0.0574869
78	32	15	6	0.0580764
79				
80	1	18	6	1.90769
81	5	18	6	0.581317
82	10	18	6	0.394283
83	15	18	6	0.305144
84	20	18	6	0.396245
85	25	18	6	0.371291
86	32	18	6	0.421195
87				
88	1	20	6	9.1156
89	5	20	6	2.6281
90	10	20	6	1.7034
91	15	20	6	1.30543
92	20	20	6	1.67951
93	25	20	6	1.53332
94	32	20	6	1.41287
95				
96	1	22	6	56.5532
97	5	22	6	13.8116
98	10	22	6	8.34245
99	15	22	6	6.20477
100	20	22	6	7.8064
101	25	22	6	6.8791
102	32	22	6	6.14886
103				
104				