A Scalable FPGA Architecture for Quantum Computing Simulation

Lee A. Belfore II

Old Dominion University, Norfolk, VA USA lbelfore@odu.edu

Abstract

A quantum computing simulation provides the opportunity to explore the behaviors of quantum circuits, study the properties of quantum gates, and develop quantum computing algorithms. Simulating quantum circuits requires geometric time and space complexities, impacting the size of the quantum circuit that can be simulated as well as the respective time required to simulate a particular circuit. Applying the parallelism inherent in the simulation and crafting custom architectures, larger quantum circuits can be simulated. A scalable accelerator architecture is proposed to provide a high performance, highly parallel, accelerator. Among the challenges of creating a scalable architecture is managing parallelism, efficiently routing quantum state components for gate evaluation, and measurement. An example is demonstrated on an Intel Agilex field programmable gate array (FPGA).

Keywords

Quantum computing, simulation, VHDL, Field programmable gate array, Scalability

Introduction

Quantum computing continues to grow in interest because of its theoretical ability to solve intractably hard computational problems that are beyond the capabilities of conventional computing methods Shor (1994); Grover (1996). Concurrently, research and industry efforts have succeeded in demonstrating increasingly capable quantum computing platforms D-Wave Systems (2020); IBM (2019b); Arute et al. (2019). Of note, quantum supremacy Preskill (2012) is the point where the quantum computing platform is able to solve problems that classical systems are not capable. Presently, the reported largest quantum computer implementations vary according to the platform & inherent capabilities and continues to increase as technological hurdles are crossed (for example IBM (2019a); Arute et al. (2019); Cao et al. (2023)) and are at the threshold of quantum supremacy.

Because access to quantum computing platforms have limited or restricted access, quantum computing simulation provides a means for people to learn about quantum computing and test their ideas within the limits of the simulator's capabilities. Further, algorithms and other ideas that may require features not present in existing quantum computing platforms, for example a novel quantum gate that is an important algorithmic building block.

Many capable quantum computing simulators are available Gheorghiu (2018); Jones et al. (2019); Xu et al. (2023). QuEST Jones et al. (2019), like other simulation platforms, includes the ability to parallelize simulations, utilizing CPU cores and GPUs, for increasingly larger problem sizes. QuEST was demonstrated to be capable of simulating quantum circuits with 38 qubits on a 2,048 node supercomputer.

Theoretically, in order to simulate any general quantum circuit, the simulator must include the ability to include gates that form a universal set Gottesman (1998); Boykin et al. (1999). As noted in Gottesman (1998), gates within the Clifford group are not fully universal and must be augmented with additional gates to make the set universal. Of importance in this work, Boykin et al. (1999) describes a set of universal gates consisting of one and two input gates.

The approach described here presents a demonstration of a quantum computing simulator implemented on a high performance field programmable gate array (FPGA) platform. FPGAs offer flexible programming available in a large and capable programmable fabric that offers the equivalent of several million logic functions. In addition, digital signal processing (DSP) units provide optimized computational modules in capacities of thousands and more. Each DSP is potentially capable of completing two multiplications every clock cycle at clock frequencies exceeding 600MHz Childs (2023). Further, block random access memory (BRAM) provides fast, flexible memory that can feed DSPs with data at high rates. Finally, FPGAs, particularly high end devices, offer high speed networking and PCIe interfaces to support distributed applications.

Building on these raw resources, FPGAs provide a flexible platform to construct arbitrary computational structures that can take leverage the unique features of a computational problem that conventional CPU & GPU based platforms cannot. The parallelism in the problem can be more directly matched to the computational problem as well as the flow of data and computational products. For example, bandwidth issues related to memory hierarchies can be minimized in FPGA implementations by placing memory containing data operands close to the computational functional units.

Several approaches for simulating quantum computers using FPGAs are described in the literature. Early results are limited as a result of the technology of the time. An early result, Khalid et al. Khalid et al. (2004) describes an approach for specifying and synthesizing quantum circuit simulators. Khalid's approach is flexible in that an arbitrary circuit can be specified, within the capabilities of the target FPGA. A weakness in Khalid's approach is that simulating a new circuit will require specifying the desired circuit and then synthesizing which can be time consuming, particularly for larger circuits. Frank et al. Frank et al. (2009) describes an approach for implementing an FPGA based quantum computing simulator but did not provide an actual demonstration implementation. Similar to Khalid, Silva and Zabeleta Silva and Zabaleta (2017) similarly synthesizes a specific circuit, the quantum Fourier transform (OFT) which is an important building block in many important quantum algorithms. Of note, the authors use Xilinx High Level Synthesis (HLS) Advanced Micro Devices (2023) to specify the circuit that is subsequently synthesized. As with Khalid et al. (2004), each new circuit must be synthesized which, for large circuits, may incur a significant amount of time for synthesis. Pilch and Długopolski Pilch and Długopolski (2019) propose a general FPGA based emulator that includes the marrying of quantum and classical procedural computing. Furthermore, their approach supports general two-input quantum gates and includes a demonstration of the Deutsch gate Deutsch (1989). While the analysis provided in the paper is comprehensive, the demonstration was for a two-qubit circuit example. Suzuki et al. (2023) presents the application of a special purpose hybrid CPU-FPGA quantum computing simulator that simulates quantum machine learning using the quantum kernel method for image classification. Notably, the simulation demonstrated the use of a 6-bit quantum kernel applied to a quantum support vector machine (QSVM).

In this paper, we present a scalable FPGA based quantum computing simulator architecture. The architecture supports any one-input quantum gate, and two-input quantum gates where the magnitude of matrix elements is one. The gates supported are consistent with universal quantum computing. The simulator is able to take an arbitrary circuit specification within the noted gate restrictions. The simulator has been implemented on an Intel Agilex FPGA capable of modeling arbitrary seven-bit quantum circuits. The results are functionally validated by comparing against the QuEST quantum logic simulator. Attempt to directly leverage FPGA resources, parallelize

This paper is organized into seven sections including an introduction, an overview of quantum circuit simulation, a review of important architectural decisions, a detailed discussion of the quantum simulation unit (QSU), a presentation of examples, and a summary & future work.

Simulating Quantum Circuits

In the context of this paper, quantum circuits are simulated on classical computing platforms. In this section, aspects of quantum circuit simulation important to this paper are present and an in depth treatment of quantum computing simulations can be found elsewhere Jones et al. (2019);

Corporation (2024). The relevant parameters, operations, and essential properties are described in this section.

Quantum State

The representation of the quantum state is central to the simulation of quantum circuits. The state of an n-bit quantum circuit is represented by a vector of 2^n probability amplitudes, where the i^{th} component holds the probability amplitude for that state. This quantum state vector (QSV) is capable of representing an arbitrary entangled state. Apropos, the QSV for an n-qubit circuit gives a snapshot of the state of the quantum state at some point in time. The QSV is expressed as

$$|\Psi\rangle = \sum_{i=0}^{2^{n-1}} \alpha_i |b_i\rangle \tag{1}$$

where n is the number of qubits, b_i is the n-bit binary code for state i, $|b_i\rangle$ is the state, and α_i is the ith complex probability amplitude. Note the geometric relationship between n and the number of components in the quantum state. Further, the norm of the state vector must be unity, or

$$\left| |\Psi\rangle \right| = \sum_{s=0}^{2^n - 1} |\alpha_i|^2 = 1.$$
 (2)

Simple Introduction to Quantum Circuit Operation

Here, a simple introduction to quantum circuit operation on one qubit is presented. The interested reader can consult Nielsen and Chuang (2010) for further details. In classical computing, information is represented as bits taking on the values '0' and '1'. Quantum information, on the other hand, is represented by qubits that are entangled mixtures of '0's and '1's. Individual qubits are represented by a one-dimensional array with two elements giving the entangled contributions for each of the binary values. The "ket" notation is used to describe the states.

$$|0\rangle = \begin{bmatrix} 1\\0 \end{bmatrix} \qquad |1\rangle = \begin{bmatrix} 0\\1 \end{bmatrix}. \tag{3}$$

Equation (3) describes what are termed "sharp" values. These are values that are known with certainty and contain no entanglements with other values or qubits.

More generally, a qubit can be described by the following equation

$$|\Psi\rangle = \alpha |0\rangle + \beta |1\rangle, \tag{4}$$

where α and β are complex values such that $|\alpha|^2 + |\beta|^2 = 1$. Equation (4) represents a value is neither '0' or '1' and that it is a quantum entanglement of these values. The entanglement can be extended to multiple bits. The following state gives the general form for the entangled state of two qubits

$$|\Psi\rangle = \alpha |00\rangle + \beta |01\rangle + \gamma |10\rangle + \delta |11\rangle,$$
 (5)

where similarly $|\alpha|^2 + |\beta|^2 + |\gamma|^2 + |\delta|^2 = 1$.

A general one qubit gate operation can be defined as an arbitrary rotation on the Bloch Sphere as the unitary matrix

that follows from

$$G_1 = U = e^{j\alpha} R_z(\beta) R_y(\gamma) R_z(\delta). \tag{6}$$

In the context of the Bloch sphere, α is the rotation about the x-axis, $R_y(\cdot)$ is a rotation about the y-axis, $R_z(\cdot)$ is a rotation about the z-axis, and β , γ , & δ and the respective rotation angles Nielsen and Chuang (2010). The unitary matrix G_1 will designate a one-input quantum gate. In the general case, the matrix elements are complex. A one-qubit gate operation on a qubit is described as

$$|\Psi'\rangle = U |\Psi\rangle \tag{7}$$

where $|\Psi'\rangle$ is the updated state after performing the gate operation. Two input gates are similarly formulated.

Universal Set of Gates

The gates available in the simulator must be carefully considered so as to provide a universal set of gates. The Clifford group, consisting of the Hadamard Gate (H), the S gate (S), and and the controlled-not gate (CNOT) gate, form the basis for many useful results but is not universal Gottesman (1998). In order to meet the criteria in the previous paragraph, the most suitable universal set is meeting the criteria of one and two input gates is given by Boykin et al. (1999), which consists of the gates H, CNOT, and the T gate (S). Beyond the minimum set of gates required for universal quantum computation, additional gates provide a means for representing quantum circuits in a more compact fashion.

Evaluation of a Quantum Gate on a Classical Computer

On a classical computer, the gate applied to the 2^n length quantum state $|\Psi\rangle$ requires "touching" each element in the vector. Assuming a one-input quantum gate applied to a qubit q_i , and further assuming the qubits are permuted so that q_i is now the right-most, $|\Psi^i\rangle$. Successive pairs of elements in $|\Psi^i\rangle$ are identical in the first (n-1) qubits state and differ in the last, nominally representing the state vector subset $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$. Applying the gate matrix to this pair of states gives the updated state given the first (n-1) are constant. Repeating for all 2^{n-1} pairs from the quantum state register results in an updated state vector reflecting the application of the gate to the original state vector. Two input gates can be similarly formulated, with a permuted state $|\Psi_n^{i,j}\rangle$, with the two-input quantum gate applied to successive quartets from the permuted quantum state register. The original quantum state is restored by applying a permutation consistent with the original ordering of the qubits. Notably, provided an efficient permutation operation and sufficient parallelism is available, performing gate operations is straightforward, with the caveat that the required parallelism is $O(2^n)$. Pseudocode for evaluating gates is outlined in Algorithms 1 and 2.

A measurement gate, or M gate, is fundamentally different in that it is not reversible and results in a discrete value, either 0 or 1, for the measured qubit. The M gate also is the means for retrieving results from the quantum circuit. Measuring a qubit on a classical computer requires three steps: quantum

Symbol	Definition			
$ \Psi angle$	QSV			
$\ket{\Psi}_i$	i th component of QSV			
$ \Psi\rangle_{ij}$	QSV segment from i to j			
$ \Psi^i angle$	permuted QSV, qubit i is LSB			
$\ket{\Psi^{ij}}$	permuted QSV, qubit i is LSB, and j next			
	most			
ψ	unnormalized segment of a QSV			
$\pi(\ket{\Psi},i)$	quantum state permutation resulting from			
	moving qubit i to the LSB. Result of			
	permutation is $ \Psi^i angle$			
$\pi^{-1}(\ket{\Psi^i},i)$	-			
	vector to its original ordering			
G_i	$2^i \times 2^i$ unitary matrix defining function for			
	an i input gate			
$\{G_1,i\}$	Specification of one input gate applied to			
	qubit i			
$\{G_2, i, j\}$	Specification of two input gate applied to			
	qubits i, j			
\mathbf{C}	Quantum circuit consisting of a Sequential			
	list of one & two input gate specifications,			
	$\{G_1, i\}$ or $\{G_2, i, j\}$, that defines a quan-			
	tum circuit quantum gate			

Table 1. Notation used in algorithmic specifications

Algorithm 1: One-input quantum gate evaluation

```
\begin{array}{l} & \text{QuantumTwo} \; (|\Psi\rangle\,,G_2,i,j) \\ \hline \textbf{Inputs} : \text{Initial quantum state} \; |\Psi\rangle \\ & \text{Two-input gate} \; G_2 \\ & \text{Qubits} \; i,j \\ \hline \textbf{Output:} \; \text{Updated quantum state} \; |\Psi'\rangle \\ & |\Psi^{ij}\rangle = \pi(|\Psi\rangle\,,i,j) \\ & \textbf{for} \; c \in \{0,1,...,2^{n-2}\} \; \textbf{do} \\ & |\psi_c|\Psi^{ij}\rangle_{4c...4c+3} \\ & |\psi_c'=G_2\psi_c| \\ & |\Psi^{ij}\rangle_{4c...4c+3} = \psi_c' \\ & \textbf{end} \\ & |\Psi'\rangle = \pi^{-1}(|\Psi^{ij'}\rangle\,,i,j) \\ & \text{return} \; |\Psi'\rangle \end{array}
```

Algorithm 2: Two-input quantum gate evaluation

state. Note that the measurement gate requires examining all elements of the QSV twice. Further, summing the probabilities serializes the function of the gate, limiting the beneficial aspects of providing hardware parallelism to complete the operation. The algorithm for evaluating an M-gate is given in Algorithm 3.

```
Measure (|\Psi\rangle, i)
Inputs: Initial quantum state |\Psi\rangle
               Qubit measured i
Output: Updated quantum state |\Psi'\rangle
|\Psi^i\rangle = \pi(|\Psi\rangle, i);
P_0=0; // Initialize zero probability
for c \in \{0, 1, ..., 2^{n-1}\} do
  | P_0 = P_0 + | |\Psi^i\rangle_{2c}|^2 
end
if P_0 \neq 0 and P_0 \neq 1 then
      prn=Rand(0,1); // uniform random number in
        [0,1]
      for c \in \{0, 1, ..., 2^{n-1}\} do
            if prn < P_0 then
                   // Qubit measured as zero
                  \begin{split} |\Psi^{i'}\rangle_{2c} = &|\Psi^i\rangle_{2c}/\sqrt{P_0} \\ |\Psi^{i'}\rangle_{2c+1} = &0 \end{split}
                   // Qubit measured as one
                   P_1 = 1 - P_0
                   \begin{aligned} &|\Psi^{i'}\rangle_{2c} = 0 \\ &|\Psi^{i'}\rangle_{2c+1} = |\Psi^{i}\rangle_{2c+1} / \sqrt{P_1} \end{aligned}
      end
else
     |\Psi^i\rangle' = |\Psi^i\rangle
|\Psi'\rangle = \pi^{-1}(|\Psi^{i'}\rangle, i)
return |\Psi'\rangle
```

Algorithm 3: Measurement gate

A quantum circuit is evaluated by sequentially applying quantum logic gates to the QSV. a sequence of gate operations applied to qubits. In the context of simulation, this naturally means that the gate operations are applied to the QSV as described in the previous section. Algorithm 4 gives the process for evaluating a quantum circuit using the previously defined algorithms. Taken together, Algorithms 1-4 define the process for simulating a quantum circuit on a classical computer.

General Architectural Decisions

In order to accelerate the simulation process, parallelism of logic and computational resources will be applied. In addition, rather than creating a general computational engine, a custom architecture is constructed specifically optimized to simulate quantum circuits. The discussions are focused more directly on midrange and higher FPGAs.

FPGA resources

FPGAs provide programmable logic, "fabric" that can be used to implement logic, switching, and state machines.

```
EvaluateCircuit (|\Psi\rangle, i)
Inputs: Initial quantum state |\Psi\rangle
           Quantum circuit C
Output: Updated quantum state |\Psi'\rangle
// assume list of quantum gates is a vector of length C
|\Psi'\rangle = |\Psi\rangle
for g \in \{0, 1, ..., C-1\} do
    if C[g] gate is a one-input gate then
         // a one-input gate
         if C[g] gate is an M-Gate then
             // apply Algorithm 3
             Measure(|\Psi'\rangle, \mathbb{C}[g].i)
         else
             // apply Algorithm 1
             QuantumOne(|\Psi'\rangle, C[g].gate, C[g].i)
         end
    else
         // otherwise, a two-input gate
         // apply Algorithm 2
         QuantumTwo(|\Psi'\rangle,C[g].gate,C[g].i,C[g].j)
    end
end
return |\Psi'\rangle
```

Algorithm 4: Evaluate quantum circuit

The fabric provides programmable units that include logic and flip-flops. In order to support high performance data processing, digital signal processors (DSPs) provide optimized & flexible multipliers, adders, registers, that provide high speed operation and efficient pipelining. A goal in FPGA and custom design is to have fast small memories, BRAM, in the vicinity of DSPs to reduce latencies incurred by sourcing operands. In order to coordinate high level operation & tasking, many FPGAs offer hard processor cores that are typically ARM based. Other specialized hard capabilities are available that provide high performance in a variety of ways including data throughput (PCIe, Ethernet), specialized AI engines, high bandwidth memories (HBMs).

Algorithm implementation on FPGAs

Algorithms on FPGAs can be implemented by mapping an algorithm to a state machine that is implemented directly on the FPGA. The state machines can be organized hierarchically in order to manage high level states related to phases of the computation and low level state machines that manage the state and computation on a per clock cycle basis.

Permutation network

From Algorithms 1-3, permutation of the inputs to the functional units is required. Notably, all algorithms perform an initial permutation to align the QSR elements to the functional units and then an inverse permutation to restore the initial order. Because the permutation is common to these algorithms, the permutation operation can be factored out, provided the permutation mechanism is general. In addition, the observant reader can note that the restoring permutation is not required provided the original permutation is remembered and each subsequent permutation is performed based on the remembered permutation. An architecture

having only one permutation network provides multiple benefits, included reducing the time necessary to route operands to functional units and reducing resources because the routing network can be expensive.

On a classical computer, the vector needs to be accessed linearly and retrieving and operating on any component requires $O(2^n)$ time. Applying a hardware solution in the form of a permutation network can take the state vector as the input and after passing through some number of stages in the permutation network, all operands for a specific gate are grouped together. Non-blocking permutation networks requiring O(logn) switching layers provide an effective solution Beneš (1964); Nassimi and Sahni (1980). Per the discussion in the previous paragraph, one permutation network is required.

One and two input gates

The simulator will model one and two input gates that form a set of gates supporting universal quantum computation. Gates beyond this set are provided for convenience to simplify the composition of quantum circuits. In order to support acceleration through parallelism, pools of one and two gates will be available that can operate in parallel. Furthermore, computational resources, i.e. addition and multiplication will be restricted to one-input gates. As noted previously, two-input gates are restricted to members of the second order Pauli group which require neither addition nor multiplication.

Result reporting module

Not shown in Figure 1 is the result reporting module (RRM). Once all qubits in the circuit are measured, the QSR will have exactly one component that is non-zero and whose magnitude is one. The RRM will identify the index associated with that component from the 2^n components in the QSR. In addition, if the QSR is not in a suitable state, i.e. two or more states are entangled, the mechanism will indicate it is incapable of providing a single result.

Numerical precision

Numerical precision can present challenges when small value differences are important or a calculation is a result of many intermediate steps. In an FPGA, the DSP provides optimized hard resources to perform addition and multiplication. In addition, the multipliers are reconfigurable so that a DSP can be reorganized into one, two, or three multipliers of varying sizes. In this demonstration, fixed point arithmetic is performed with an integer precision of one bit and fractional precision of 16 bits for multiplication operations. With this precision, a typical DSP can provide two independent multiplications. For situations where a result is a combination of many operands, extended precision is provided, for example during measurement operations.

Verification

Verifying the operation of the quantum computing simulation accelerator will be accomplished by comparing results with an established quantum computing simulation platform, such as QuEST Jones et al. (2019).

Quantum Simulation Unit Architecture

In this section, an architecture, suitable for implementation on an FPGA or custom silicon platform is presented. The Figure 1 gives the architecture that follows from the discussions in the previous two sections. Omitted from this figure is an embedded processor that is used to control & monitor the operation of the QSU. The demonstration has been implemented in an Intel Agilex 7 FPGA which includes a four core ARM hard processor system (HPS).

With trivial modifications, all constant one-input gates can be easily added and all two-input gates where the magnitude of matrix coefficients is one can be added.

An embedded processor controls the operation of the QSU implemented on the FPGA fabric through an interface that consists of four 32-bit registers. The registers and their purposes are described in Table 2. In principal, the embedded process could be a soft processor core or a hard processor system (HPS). In this work, an HPS running the Yocto Linux operating system fills this role. The HPS_control register is a write only and controls the overall operation of the QSU, sets the mode, configures the QSU to receive input, and report outputs. The HPS_status register is read only and is used to monitor the operation of the QSU. The HPS_data and Fabric_data pass data to and from the QSU. For example, the list of simulated gates is passed from the HPS to the fabric through HPS_data. Simulated results and debugging data can be sent from the QSU to the HPS through Fabric data.

Register name	Mode	Purpose	
HPS_control	write only	control over fabric	
Fabric_status	read only	status from fabric	
HPS_data	write only	supply data to fabric	
Fabric_data	read only	receive data from fab-	
		ric	

Table 2. HPS QSU interface

The QSU architecture follows from Algorithms 1-4 with the overarching goal to leverage their inherent parallelism. The architecture is organized into several modules and its datapath is shown in Figure 1. To the extent possible, modules are designed to operate independently from the other modules and recognize when operands and other relevant information are available to proceed with their inherent function. The control path consists of a two level hierarchy of state machines that controls sequencing of data through the datapath. The top level state machine follows from Algorithm 4 which controls state machines for Algorithms 1-3. Because the permutation is common to both one and two inputs, that function has been factored out and is controlled by the top level state machine.

Quantum State Register (QSR)

The QSR holds the QSV for the n qubit circuit. It is an array of 2^n complex fixed point values. The QSR reflects the original ordering of qubits and not any permutation. Because many existing FPGA platforms do not directly support floating point and to offer better utilization of DSPs, fixed point arithmetic is used.

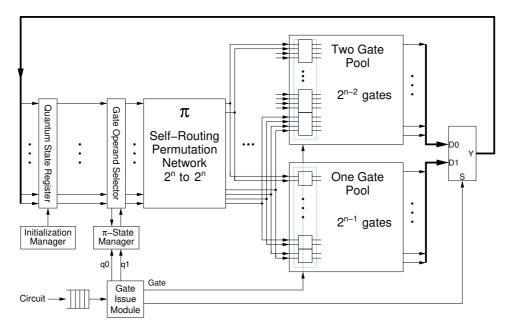


Figure 1. Top level quantum simulation unit (QSU)

Initialization Manager (IM)

The initialization manager provides the capability to an arbitrary initialization of the QSR to an arbitrary entangled initial state. In the event the initialization provided is unnormalized, the IM will initiate normalize the state by initiating a special normalization operation within the One-Qubit Gate Pool module.

Gate Issue Module (GIM)

The GIM consists of two parts. The first part is a buffer and associated buffer management that stores the list if gates to be evaluated. The buffer is loaded through the HPS fabric interface and new gates to be evaluated can be added at any time. The gate issued is at the head of the buffer and is partially decoded to retrieve the gate to be evaluated as well as its required operands. Also included is the implementation of Algorithm 4.

Gate Operand Selector (GOS)

As required in Algorithms 1 and 2, permutations are required to forward the appropriate components from the QSR to the functional units that implement gate functions. The GOS packages the components with current gate information, qubit inputs, and final permutation position necessary for the permutation.

Self-Routing Permutation Network (SRPN)

As noted, the permutation network must forward the operands to the appropriate gate functional unit. Because this requires a full 2^n to 2^n permutation, an efficient permutation implementation is required. The fastest permutation network is a crossbar switch. While the crossbar switch can perform the permutation in one step, the number of switches required is 2^{2n} which is undesirable. A reduction in the number of switches is possible using a switching network. While requiring more time to perform the permutation, the number of switches is $O(2^n)$. A further challenge in

performing the permutation is "collisions" that result in blocking can delay passage through the switching network. A well known non-blocking switching network is the Beneš network Beneš (1964) which enables arbitrary permutations using (n/2-1) layers of 2^{n-1} switching elements. The functionality required in the switches is described elsewhere Nassimi and Sahni (1980). The reader may note that the permutation required here is not totally general and has structure that may reduce the number of layers and will be addressed in future work. While Algorithms 1, 2, & 3 include permutation and inverse permutation operations, Figure 1 only shows the inclusion of one permutation network. One permutation network is sufficient provided the permutation is known when each gate is simulated and the original permutation is restored after the circuit has been simulated. An example permutation network is shown for a 2^4 to 2^4 instance in Figure 2. Because arbitrary sized Beneš networks can be defined recursively, generalized HDL specifications are possible. Furthermore, the Beneš network is capable of partial permutations, i.e. 2^i to 2^j where j < i in the event resources are limited permitting pipelined functional units for faster processing.

One-Qubit Gate Pool (1-QGP)

The 1-QGP holds the one-input gate functional units and controls the operation of the gates within the pool. In order to accelerate the simulation, gates within the pool operate in parallel on the permuted QSR. Figure 3 gives the architecture for the gate pool. For non-measurement gates, the results are passed on to the network that performs the inverse permutation. For measurement gates, applying the gate is more complex, requiring several global operations that must be performed. First, per Algorithm 3, the probability that the measurement results in 0 is calculated. Notably, measurement gate implementation challenges were observed by others Pilch and Długopolski (2019). The operation requires computing the magnitudes of the QSR components which is performed in the one-input gate functional unites,

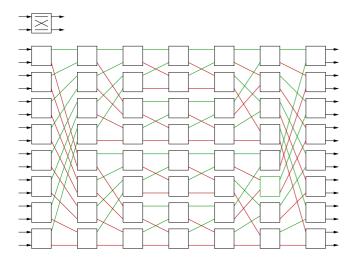


Figure 2. Self-routing permutation network for n=4

summing the values with a summing network within the 1-QGP. Once the probability is computed, several special cases need to be handled. If the measured gate already has a sharp value, no assignment of the qubit is necessary. If the probability is neither 0 or 1, the qubit is collapsed by comparing a pseudorandom number with the probability. The collapse requires that 2^{n-1} components be zeroed as a result of one qubit value now being excluded from possibility. At this point, the QSR is no longer properly normalized, and the square root of the winning bit's probability must be computed followed by computing its inverse.

Because of the likelihood of numerical errors, computed probabilities may not sum exactly to precise values as implied by the mathematics. This situation is mitigated, but not entirely eliminated, by using extended precision arithmetic in the adder network, doubling the precision of the fixed point number and including an error tolerance that is a function of the number of qubits in the circuit. For the inverse calculation, because the inverse is greater than 1, an extended precision type capable of holding these values is used.

As noted for universal computation, in support of the minimum set of quantum gates required, the one-input Hadamard and T gates are required. Further, the measurement or M gate is required to retrieve answers. For convenience to support a wide variety of quantum circuits, many more gates have been implemented and are listed in Table 3 where $j = \sqrt{-1}$.

The required arithmetic operations local to each gate are addition and multiplication. Rotation gates are common in many quantum algorithms. Any constant rotation can be included by specifying its unitary matrix. At the present, the simulation is not capable of handling an arbitrary matrix. If this is desired, a simple fix would be to augment the HPS-fabric interface to permit the HPS to provide any desired unitary matrix.

Two-Qubit Gate Pool (2-QGP)

The 2-QGP models two-input gates in a similar fashion to one-input gates. The permutation network routing differs because the two-input gates require four components from the QSR. As noted, multiplications are restricted to the

Gate	Name	Function
Nop	No-op gate	$\left[\begin{array}{cc} 1 & 0 \\ 0 & 1 \end{array}\right]$
X	Pauli X	$\left[\begin{array}{cc} 0 & 1 \\ 1 & 0 \end{array}\right]$
Y	Pauli Y	$\left[\begin{array}{cc}0&j\\-j&0\end{array}\right]$
Z	Pauli Z	$\left[\begin{array}{cc} 1 & 0 \\ 0 & -1 \end{array}\right]$
Н	Hadamard	$\frac{1}{\sqrt{2}} \left[\begin{array}{cc} 1 & 1\\ 1 & -1 \end{array} \right]$
V	$\sqrt{\text{NOT}}, \sqrt{X}$	$\frac{1}{2} \left[\begin{array}{cc} (1+j) & (1-j) \\ (1-j) & (1+j) \end{array} \right]$
\sqrt{Y}		$\frac{1}{2} \left[\begin{array}{cc} (1+j) & (-1-j) \\ (1+j) & (1+j) \end{array} \right]$
S	\sqrt{Z}	$\left[\begin{array}{cc} 1 & 0 \\ 0 & j \end{array}\right]$
S^{-1}	Inverse of S	$\left[\begin{array}{cc} 1 & 0 \\ 0 & -j \end{array}\right]$
T	\sqrt{S}	$\left[\begin{array}{cc} 1 & 0 \\ 0 & e^{j\pi/4} \end{array}\right]$
T^{-1}	Inverse of T	$\left[\begin{array}{cc} 1 & 0 \\ 0 & e^{-j\pi/4} \end{array}\right]$
E_x	Error in X	Apply Pauli X with a designated error probability
E_y	Error in X	Apply Pauli Y with a designated error probability
E_z	Error in Z	Apply Pauli Z with a designated error probability
M	Measurement gate	Measures and collapses the state for the specified qubit

Table 3. Implemented one-input quantum gates

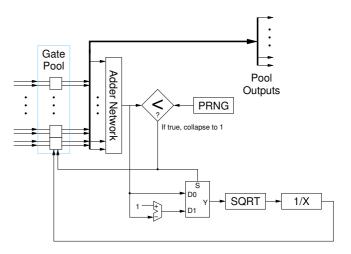


Figure 3. One qubit gate pool block diagram

one-input gates, limiting the nature of two-input gates. The unitary matrices defining two-input gates requires matrix elements be restricted to $\{1,-1,j,-j\}$. A summary of the two-input gates implemented are shown in Table 4. Overall, the architecture of 2-QGP is relatively simple compared with the 1-QGP.

Universal Set of Gates

A review of Tables 3 and 4 reveals that the QSU includes gates that support a universal set Boykin et al. (1999) since

Gate	Name	Function
CNOT	Controlled-NOT	$ \left[\begin{array}{ccccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{array}\right] $
CY	Controlled Pauli Y	$\left[\begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & j \\ 0 & 0 & -j & 0 \end{array}\right]$
CZ	Controlled Pauli Z	$ \left[\begin{array}{ccccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 \end{array}\right] $
\sqrt{ZZ}	$\sqrt{Z \bigotimes Z}$	$ \left[\begin{array}{ccccc} 1 & 0 & 0 & 0 \\ 0 & j & 0 & 0 \\ 0 & 0 & j & 0 \\ 0 & 0 & 0 & -1 \end{array}\right] $
SWAP	Swap gate	$\left[\begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{array}\right]$

Table 4. Two-input gates

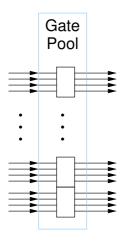


Figure 4. Two qubit gate pool block diagram

CNOT gate, Hadamard Gate and the T are included in the set of gates implemented.

Example

The QSU is implemented at the register transfer language (RTL) level using 2008 standard release of the VHSIC Hardware Description Language (VHDL) IEEE Computer Society (2009). The QSU is verified through simulation and also on an FPGA for two circuits and the results are described in this section.

Simulation Results

Simulation was performed first on the RTL model and then confirmed using the QuEST Jones et al. (2019) quantum computing simulator. The results for two circuits are given here.

Simple Circuit. VHDL simulation results have been performed using the open source VHDL platform GHDL v2.0 Lehman et al. (2022) for a simple circuit composed of a layer of Hadamard gates to fully entangle the quantum state is applied. This is followed by a measurement of one qubit, a

controlled-NOT gate, and then the subsequent measurement of the remaining qubits. The circuit is shown in Figure 5.

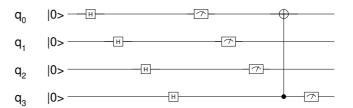


Figure 5. Test circuit schematic

Random Circuit. Random circuits provide a stronger verification of the QSU architecture's correct operation because biases are eliminated in the creation of the circuit. In addition, random quantum circuits have benefit in their own right in regards to solving certain classes of problems (see for example Emerson (2004); Fisher et al. (2023). A 7-qubit random circuit example is shown in Figure 6 using an approach outlined in Emerson (2004). The circuit is generated by first applying the Hadamard gate to each input, fully entangling all inputs. Next, for three iterations, CNOT gates are applied as shown followed by, for each qubit, gates randomly selected from the list $\{X, Y, Z, S, S^{-1}, T, T^{-1}, \sqrt{X}, \sqrt{Y}, \sqrt{Z}\}$. In Figure 6, the staggered order of the gates reflects their simulation order. Once the circuit has been defined, it is included in a test bench for simulation of the QSU HDL specification. Aggregate statistics are generated regarding the final output of the circuit for ten thousand trials. Next, the simulation is performed using QuEST. More simulation trials, one million, are run to achieve a more refined estimate of the expected distribution among final states. Finally, the QSU is programmed onto the FPGA and the circuit is similarly run for one million trials. The results are summarized in Table 5.

Qualitatively comparing the QuEST and QSU FPGA based results, the respective empirically derived state probabilities are very similar and are zero for the same quantum states. Further, the Euclidean distance between the distributions is 1.45×10^{-4} , mean absolute error between the distributions is 1.36×10^{-5} with a standard deviation of 1.21×10^{-5} . These suggest very good agreement between the two simulation results.

FPGA Demonstration

The QSU has been implemented and evaluated on an Intel®Agilex®AGFB014R24A2E2VR0 (AGFB014) FPGA on the Intel®Agilex®F-Series FPGA development kit Intel (2023a). The FPGA programming has been developed using the Intel®Quartus®Prime Pro version 21.1 development platform Intel (2023b). The AGFB014 fabric includes 487,200 Adaptive Logic Modules (ALMs) (1,437,240 logic element equivalent), 145.6Mb of BRAM, 4,510 digital signal processors (DSPs). In addition, the AGFB014 includes hard IP Ethernet capable of supporting data rates up to 400Gbps data rates and PCIe 4.0. In addition, the HPS is a quad core 64-bit Arm Cortex-A53.

Table 6 provides a summary of the utilization of some major QSU modules for a QSU capable of simulating a 7-qubit circuit.

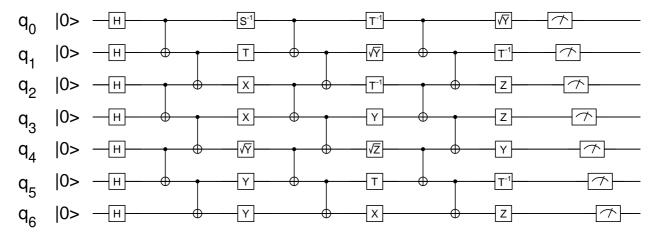


Figure 6. Random circuit

Performance

Few attempts were made to optimize the performance and space of the QSU and, rather, creating a working prototype was prioritized. Making it is challenging to make direct comparisons. Simulating the random circuit required 1,430 clocks per individual circuit simulation. With a clock period of 10ns, the time required for an individual simulation is $14.3\mu s$ Notably, simulating measurement gates consumed 271 clocks, or nearly 20% of the simulation using QuEST on a XEON processor (XEON ES-2603 v2, clock frequency of 1.8GHz) requires 5.25ms of CPU time/simulation. The FPGA based simulation provides a speedup of 368.

Summary and Future Work

In this paper, an FPGA based quantum computing simulator was presented. The simulator models sufficient gates to support universal quantum computation. Further, the simulator architecture was conceived to be extensible and scalable. The simulator was modeled at the RTL level using VHDL. A demonstration implementation on an Intel Agilex FPGA was capable of simulating a seven qubit quantum circuit. The FPGA results were verified by comparing with simulations both of the VHDL implementation and with respect to the QuEST quantum computing simulator.

Several potential directions for future work are considered. First, the model presented here has been implemented as a proof of concept and beyond architectural features & use of massive parallelism, few attempts were made to optimize performance & resources. Of course, simulating an arbitrary quantum circuit is at least linear in time & resources with respect to the length of the state vector. The time required to perform certain serial aspects, such as the square root & division required in the measurement gate can be improved. In addition, the number of DSPs required for each individual quantum gate can potentially increase the potential parallelism. Second, demonstrating scalability across multiple devices would enable simulating larger problem sizes. Third, the quantum circuit simulation can be coupled with classical computing algorithms to simulate algorithms that include both classical procedural and quantum aspects. Fourth, additional flexibility can be provided by providing a capability of loading arbitrary unitary matrices for one-input gates. For example, algorithms such as the quantum Fourier transform (QFT) Coppersmith (2002) require several different phase gates, as determined by the size of the QFT, which are not directly supported in the QSU.

References

Advanced Micro Devices (2023) Introduction to Vitis HLS. URL https://docs.xilinx.com/r/en-US/ug1399-vitis-hls/Introduction-to-Vitis-HLS. Accessed: 2023-09-10.

Arute F, Arya K and al e (2019) Quantum supremacy using a programmable superconducting processor. *Nature* 574: 505– 510.

Beneš V (1964) Optimal rearrangeable multistage connecting networks. *The Bell System Technical Journal* 43(4): 1641–1656.

Boykin P, Mor T, Pulver M, Roychowdhury V and Vatan F (1999)
On universal and fault-tolerant quantum computing: a novel basis and a new constructive proof of universality for shor's basis. In: 40th Annual Symposium on Foundations of Computer Science (Cat. No.99CB37039). pp. 486–494. DOI:10.1109/SFFCS.1999.814621.

Cao S, Wu B, Chen F and et al (2023) Generation of genuine entanglement up to 51 superconducting qubits. *Nature* 619: 619–742. DOI:10.1038/s41586-023-06195-1. Doi:10.1038/s41586-023-06195-1.

Childs N (2023) How Intel Agilex[®] FPGA is enabling resource and power efficient 4K, 8K video processing solutions. Intel. White paper

Coppersmith D (2002) An approximate fourier transform useful in quantum factoring. Doi:10.48550/arXiv.quant-ph/0201067.

Corporation IBM (2024) Qiskit. URL https://www.ibm.com/quantum/qiskit. Accessed: 2024-03-05.

D-Wave Systems (2020) Practical quantum computing: D-Wave technical overview. URL https://www.dwavesys.com/media/doaarwgt/14-1047a-a_practical_quantum_computing_an_update.pdf. Accessed: 2023-08-14.

Deutsch DE (1989) Quantum computational networks. *Proceedings of the Royal Society A* 425(1868): 73–79. DOI:10.1098/rspa.1989.0099.

State	QSU (HDL)		QuEST		QSU (FPGA)	
$ \psi\rangle$	n/10 ⁴	$\mu_{ \psi\rangle}$	n/10 ⁶	$\mu_{ \psi angle}$	n/10 ⁶	$\mu_{ \psi angle}$
0x10	281	0.00281	23738	0.00237	23420	0.00234
0x11	76	0.00076	7869	0.00079	7938	0.00079
0x12	84	0.00084	7792	0.00078	7895	0.00079
0x13	252	0.00252	23489	0.00234	23697	0.00237
0x14	86	0.00086	7797	0.00078	7888	0.00079
0x15	233	0.00233	23478	0.00234	23587	0.00236
0x16	217	0.00217	23287	0.00232	23493	0.00235
0x17	74	0.00074	7912	0.00079	7971	0.00080
0x18>	229	0.00229	23420	0.00234	23412	0.00234
0x19>	76	0.00076	7814	0.00078	7830	0.00078
0x1a>	77	0.00077	7783	0.00078	7739	0.00077
0x1b>	225	0.00225	23348	0.00233	23384	0.00234
0x1c>	85	0.00085	7690	0.00077	7880	0.00079
0x1d>	249	0.00249	23348	0.00233	23676	0.00237
0x1e	224	0.00224	23703	0.00237	23803	0.00238
0x1f>	82	0.00082	7795	0.00078	7836	0.00078
0x30	229	0.00229	23495	0.00235	23627	0.00236
0x31	87	0.00087	7755	0.00076	7677	0.00077
0x32	92	0.00092	7840	0.00078	7835	0.00078
0x33}	240	0.00240	23369	0.00234	23470	0.00235
0x34\	74 224	0.00074	7762	0.00078	7735	0.00077 0.00235
0x35}		0.00224	23227	0.00232	23466	
0x36\ 0x37\	231 100	0.00231 0.00100	23449 7790	0.00234 0.00079	23364 7836	0.00234 0.00078
0x37) 0x38	215	0.00100	23378	0.00079	23563	0.00078
0x36/ 0x39	74	0.00213	7840	0.00234	7773	0.00230
0x39/ 0x3a/	74 76	0.00074	7766	0.00078	7790	0.00078
0x3a/	253	0.00253	23312	0.00233	23374	0.00234
0x3c/	74	0.00233	7678	0.00233	7761	0.00234
0x3d	224	0.00224	23451	0.00235	23429	0.00234
0x3e	228	0.00228	23560	0.00236	23292	0.00233
0x3f	73	0.00073	7727	0.00077	7989	0.00080
0x40	224	0.00224	23380	0.00234	23484	0.00235
0x41	81	0.00081	7840	0.00078	7878	0.00079
0x42	76	0.00076	7837	0.00078	7868	0.00079
0x43	238	0.00238	23427	0.00234	23376	0.00234
0x44>	82	0.00082	7829	0.00078	7757	0.00078
0x45}	228	0.00228	23513	0.00235	23643	0.00236
0x46}	227	0.00227	23325	0.00233	23519	0.00235
0x47>	83	0.00083	7801	0.00078	7797	0.00078
0x48}	247	0.00247	23621	0.00236	23460	0.00235
0x49>	75	0.00075	7867	0.00079	7678	0.00077
0x4a>	86	0.00086	7805	0.00078	7723	0.00077
0x4b	232	0.00232	23349	0.00233	23295	0.00233
0x4c>	54	0.00054	7802	0.00078	7780	0.00078
0x4d	215	0.00215	23661	0.00237	23438	0.00234
0x4e	248	0.00248	23174	0.00232	23330	0.00233
0x4f\	68 246	0.00068 0.00246	7835 23137	0.00078 0.00231	7704 23451	0.00077 0.00235
0x60} 0x61}	246 75	0.00246	7736	0.00231	7885	0.00233
0x61) 0x62)	96	0.00073	7851	0.00077	7748	0.00079
0x62) 0x63	242	0.00090	23271	0.00079	23320	0.00077
0x63/	84	0.00242	7874	0.00233	7718	0.00233
0x65	265	0.00265	23448	0.00234	23398	0.00234
0x66	225	0.00225	23447	0.00234	23228	0.00231
0x67	83	0.00083	7960	0.00080	7783	0.00078
0x68	217	0.00217	23560	0.00236	23230	0.00232
0x69	59	0.00059	7828	0.00078	7750	0.00078
0x6a	78	0.00078	7761	0.00078	7791	0.00078
0x6b	237	0.00237	23345	0.00233	23489	0.00235
0x6c)	67	0.00067	8057	0.00081	7741	0.00077
0x6d)	213	0.00213	23524	0.00235	23434	0.00234
0x6e	234	0.00234	23546	0.00235	23098	0.00231
0x6f>	71	0.00071	7927	0.00079	7776	0.00078
others	0	0.00000	0	0.00000	0	0.00000
						-

Table 5. Summary of verification results, rounded to five places

Model	Utilization
1-QGP	72,894 ALMs
	1,600 DSPs
2-QGP	8,333 ALMs
π	61,298 ALMs
GIM/GOS	505 ALMs
RRM	1,328 ALMs
miscellaneous	21,232 ALMs
Total	168,931 ALMs

Table 6. Module utilization

Emerson J (2004) Random quantum circuits and pseudo-random operators: Theory and applications. Doi:10.48550/arXiv.quantph/041008.

Fisher MP, Khemani V, Nahum A and Vijay S (2023) Random quantum circuits. *Annual Review of Condensed Matter Physics* 14: 335–379.

Frank MP, Oniciuc L, Meyer-Baese UH and Chiorescu I (2009)
A space-efficient quantum computer simulator suitable for high-speed FPGA implementation. In: Donkor EJ, Pirich AR and Brandt HE (eds.) SPIE Proceedings. SPIE. DOI: 10.1117/12.817924. URL https://doi.org/10.1117 %2F12.817924.

Gheorghiu V (2018) Quantum++: A modern C++ quantum computing library. *PLoS ONE* 13(12). DOI:10.1371/journal.pone.0208073. Doi:10.1371/journal.pone.0208073.

Gottesman D (1998) The Heisenberg representation of quantum computers. Doi:10.48550/arXiv.quant-ph/9807006.

Grover LK (1996) A fast quantum mechanical algorithm for database search. In: *Proceedings of the 28th annuam ACM Symposium on the Theory of Computing (STOC)*. pp. 212–219.

IBM (2019a) IBM opens quantum computation center in New York;
Brings world's largest fleet of quantum computing systems online, Unveils new 53-qubit quantum system for broad use. https://newsroom.ibm.com/2019-09-18-IBM-Opens-Quantum-Computation-Center-in-New-York-Brings-Worlds-Largest-Fleet-of-Quantum-Computing-Systems-Online-Unveils-New-53-Qubit-Quantum-System-for-Broad-Use.

IBM (2019b) IBM unveils world's first integrated quantum computing system for commercial use. URL https://newsroom.ibm.com/2019-01-08-IBM-Unveils-Worlds-First-Integrated-Quantum-Computing-System-for-Commercial-Use. Accessed: 2023-08-14.

IEEE Computer Society Design Automation Standards Committee (2009) IEEE standard VHDL language reference manual. *IEEE Std 1076-2008* DOI:10.1109/IEEESTD.2009.4772740. Doi:10.1109/IEEESTD.2009.4772740.

Intel (2023a) Intel Agilex® 7 FPGA F-Series development kit (P-tile and E-tile). https://www.intel.com/content/www/us/en/products/details/fpga/development-kits/agilex/agf014.html. Accessed: 2023-07-28.

Intel (2023b) Intel® Quartus® Prime software. https://
www.intel.com/content/www/us/en/products/
details/fpga/development-tools/quartusprime/resource.html. Accessed: 2023-07-28.

- Jones T, Brown A, Bush I and Benjamin SC (2019) QuEST and high performance simulation of quantum computers. *Scientific Reports* 9.
- Khalid AU, Zilic Z and Radecka K (2004) FPGA emulation of quantum circuits. In: *Proceedings of the IEEE International Conference on Computer Design ICCD'04*. San Jose, CA, pp. 310–315.
- Lehman P, Gingold T and Martinez-Corral U (2022) GHDL release 2.0.0. https://github.com/ghdl/ghdl/releases/tag/v2.0.0.
- Nassimi D and Sahni S (1980) A self routing Benes network. In: *Proceedings of the 7th annual symposium on Computer Architecture*. pp. 190–195.
- Nielsen MA and Chuang IL (2010) *Quantum Computation and Quantum Information: 10th Anniversary Edition.* Cambridge University Press.
- Pilch J and Długopolski J (2019) An FPGA-based real quantum computer emulator. *Journal of Computational Electronics* 18(1): 329–342.

- Preskill J (2012) Quantum computing and the entanglement frontier. Doi:10.48550/arXiv.1203.5813.
- Shor P (1994) Algorithms for quantum computation: discrete logarithms and factoring. In: *Proceedings 35th Annual Symposium on Foundations of Computer Science*. Santa Fe, NM, pp. 124–134. DOI:10.1109/SFCS.1994.365700. Doi:10.1109/SFCS.1994.365700.
- Silva A and Zabaleta OG (2017) FPGA quantum computing emulator using high level design tools. In: 2017 Eight Argentine Symposium and Conference on Embedded Systems (CASE). pp. 1–6. DOI:10.23919/SASE-CASE.2017.8115369. Doi:10.23919/SASE-CASE.2017.8115369.
- Suzuki T, Miyazaki T, Inaritai T and Otsuka T (2023) Quantum AI simulator using a hybrid CPU–FPGA approach. *Scientific Reports* 13. DOI:10.1038/s41598-023-34600-2.
- Xu X, Benjamin S, Sun J, Yuan X and Zhang P (2023) A herculean task: Classical simulation of quantum computers. Doi:10.48550/arXiv.2302.08880.