

Improvements on the Efficiency of the Practical Execution of Hybrid Quantum-Classical Algorithms

Thesis submitted by

Steven Oud

under the guidance of

**Harold Meerwaldt, QuTech
Damian Podareanu, SURF
Ariana Torres Knoop, SURF
Richard Versluis, TNO**

in partial fulfillment of the requirements for the degree of

Bachelor of Science

Improvements on the Efficiency of the Practical Execution of Hybrid Quantum-Classical Algorithms

Steven Oud*
500776959

*Faculty of Computer Science, Information Technology,
Business IT and Management
Software Engineering*

Advisor: Marten Teitsma

Amsterdam University of Applied Sciences
May 31, 2021

Abstract

Quantum computers promise to solve certain problems more efficiently than classical computers by making use of quantum mechanical phenomena. However, current quantum computers are small in qubit count and are prone to errors, greatly limiting the number of useful applications. In term, hybrid quantum-classical algorithms (HQCAs) that use classical optimization to train a parameterized quantum circuit are being researched and developed. Hybrid quantum-classical algorithms have been proposed for most applications envisioned for quantum computers, and they are a promising candidate for being the first practical application of quantum computers. In this report, we analyze the efficiency of different workflows for the practical execution of HQCAs using both quantum simulation and quantum processing units. We recommend different methods for improving the efficiency of the execution of HQCAs using the Quantum Inspire quantum computing platform and SURF's high performance computing center. While these recommendations are made with these specific platforms in mind, most of it is applicable to any hybrid quantum-classical setting.

Contents

List of Figures	iv
List of Tables	v
1 Introduction	1
1.1 Motivation	1
1.2 Objective	2
1.3 Outline	2
2 Background	4
2.1 Computational Complexity	4
2.1.1 Big-O Notation	4
2.1.2 Turing Machines	5
2.1.3 Complexity Classes	6
2.2 Quantum Mechanics	8
2.2.1 Qubits	8
2.2.2 State Evolution	10
2.2.3 Measurement	11
2.3 Quantum Computation	12
2.3.1 Quantum Circuits	13
2.3.2 Universal Gate Sets	15
2.3.3 Quantum Algorithms	15
3 Hybrid Quantum-Classical Algorithms	20
3.1 Variational Hybrid Quantum-Classical Algorithms	20
3.1.1 Cost Function	21
3.1.2 Ansätze	21
3.1.3 Optimization	22
3.2 Variational Quantum Eigensolver	23
3.3 Quantum Approximate Optimization Algorithm	25
4 Practical Hybrid Quantum-Classical Computing	27
4.1 Quantum Inspire	27
4.2 SURF High Performance Computing	29
4.3 Hybrid Quantum-Classical Workflows	29
4.4 Hybrid Quantum-Classical Workflow Analysis	30
4.4.1 Development Cycle	32
4.4.2 QX Simulator Batch Parameters	32
4.4.3 Parallel Quantum Circuit Execution	37
4.4.4 Quantum Circuit Hardware Batching	38

4.4.5 Summary	38
5 Conclusion	40
Bibliography	41
Glossary	46
Appendices	47
A QAOA for Max-Cut	48
A.1 Implementation	49
A.2 Results	51

List of Figures

1.1	General structure of a HQCA.	2
2.1	The tape of a single-tape Turing machine in an arbitrary state.	6
2.2	An overview of the hierarchy of some of the complexity classes discussed. . .	7
2.3	Bloch sphere representation of a qubit's state.	9
2.4	Quantum circuit for creating and measuring a Bell state $ \Phi^+\rangle$	14
2.5	A quantum circuit consisting of 8 gates with a depth of 4.	14
2.6	Copying a computational basis state using CNOT where $j \in \{0, 1\}$	15
2.7	Decomposition of the Toffoli gate using CNOT, H , S and T gates.	15
2.8	General circuit for the QFT on 2^n amplitudes with n qubits.	16
2.9	Quantum circuit representation of Grover's algorithm.	18
3.1	The general structure of a variational HQCA.	21
3.2	Visualization of gradient descent finding a local minimum of a function. . . .	23
4.1	Overview of the Quantum Inspire workflow.	28
4.2	Visualizations of the QX simulator benchmark results for $n \geq 30$	35
4.3	Plot of the RSS memory usage of QX simulator as the number of qubits increases. 37	
A.1	Max-Cut example on a graph with five vertices and unit weights.	48
A.2	Undirected 2-regular graph $G = (V, E)$ with $n = 4$ vertices $V = \{0, 1, 2, 3\}$ and 4 edges $E = \{(0, 1), (0, 3), (1, 2), (2, 3)\}$ with unit weight $w_{j,k} = w_{k,j} = 1$	49
A.3	Quantum circuit for the p -layer QAOA on the graph from Figure A.2.	50
A.4	cQASM for the p -layer QAOA on the graph from Figure A.2.	50
A.5	Experimental results of the running the QAOA to solve Max-Cut problem on the graph from Figure A.2.	52

List of Tables

2.1	Common big-O run times from fast to slow.	5
2.2	Frequently used quantum gates with their circuit symbol and matrix representation.	13
4.1	Specification of the Lisa and Cartesius cluster computers.	29
4.2	Relevant hybrid quantum-classical workflows and their use cases.	30
4.3	Benchmarks of the time taken by different quantum parts in the context of a hybrid quantum-classical computation.	31
4.4	Overview of the layout of the sockets and CPUs on the <code>fat_soil_shared</code> quantum HPC nodes.	33
4.5	Benchmark results of the execution time of quantum circuits using QX simulator for different qubit and core counts.	34
4.6	QX simulator RSS memory usage for varying number of qubits.	36

1

Introduction

As an introduction to the work that was done during this thesis, this chapter starts with presenting a brief overview on the relevance of quantum computing, and describing the motivation behind this work. After that, the objective of this work is described and an outline of this report is given.

1.1 Motivation

Quantum computing is a recently proposed model of computation that promises to solve certain problems more efficiently than classical computers by making use of quantum mechanical phenomena such as superposition, entanglement, and interference. The idea of quantum computing originated from Benioff [1], who proposed a quantum mechanical model of the Turing machine in 1980. This idea was later extended as Manin [2] and Feynman [3] independently suggested that quantum computers have the potential to solve certain computational problems intractable by classical computers. Since then, researchers have been searching for applications for quantum computing. Some noteworthy developments in the field of quantum computing include Shor's algorithm for factoring integers [4] and Grover's algorithm for unstructured database search [5]. These quantum algorithms promise an exponential and quadratic speedup respectively over their best-known classical counterparts. The finding of such speedups have catalyzed research towards quantum computers, and more applications have since been found in fields including chemistry [6], cryptography [7], and machine learning [8].

Until recently, quantum computing had been a mainly theoretical field. These days however, thanks to recent technological advances, various quantum devices are being actively developed. Furthermore, technological giants such as IBM, Microsoft, Intel, and Google are investing heavily in the development of quantum computers. The applications of these current quantum devices are very limited however, due to their inability to store information for a long time and their sensitivity to errors. For instance, Fowler, Mariantoni, Martinis, *et al.* [9, Appendix M] estimate that factoring a 2000-bit number using Shor's algorithm would require a quantum computer with about a billion physical qubits, and error rates below 4×10^{-13} . In contrast, devices today have about 50 to 100 physical qubits with error rates above 0.1.

These small and error-prone quantum devices, often referred to as noisy intermediate-scale quantum (NISQ) devices, may be unfit to run quantum algorithms like Shor's

algorithm, but they may still prove to be useful and perform tasks intractable by classical computers [10, 11]. To deal with the limitations of NISQ devices, hybrid quantum-classical algorithms (HQCAs) are being actively researched. These HQCAs combine classical and quantum computations as visualized in Figure 1.1. Hybrid quantum-

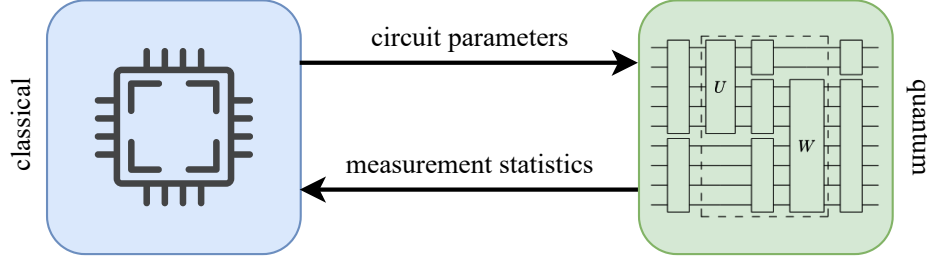


Figure 1.1: General structure of a HQCA. The data interchange between the classical and quantum part is repeated many times.

classical algorithms typically involve a small quantum computation inside of a classical optimization loop, greatly reducing the amount of quantum resources needed. This makes them suitable to run on NISQ devices, and are expected to be one of the first useful applications for quantum computing [12]. It is important to note that NISQ devices running HQCAs will likely not be revolutionary by itself. Instead, it should be seen as an important stepping stone towards more powerful quantum devices and algorithms.

Research towards HQCAs often involves executing quantum circuits on actual quantum processing units (QPUs) and through quantum circuit simulation. To support this kind of research, classical and quantum computing facilities are needed. Furthermore, these facilities need to be connected and able to interchange data in a timely manner for this kind of research to be feasible. This is especially challenging given that both quantum and classical resources are shared with other users.

1.2 Objective

The purpose of this thesis is to analyze different workflows for the practical execution of HQCAs and propose improvements to be made to increase the efficiency. This work is in collaboration with TNO, QuTech, and SURF, and is focused on the efficient execution of HQCAs using QuTech’s quantum computing platform Quantum Inspire [13] and SURF’s high performance computing (HPC) center [14]. To allow for the efficient execution of HQCAs, the SURF HPC and Quantum Inspire computing facilities should be able to interchange data in a timely manner. A key part in this is figuring out sources of overhead and current bottlenecks. In a hybrid setup like this, overheads such as quantum and classical scheduling wait times, data transfers, and resource initialization can quickly increase the run time of HQCAs.

1.3 Outline

This report is structured as follows. Chapter 2 provides the reader with a background on computational complexity, quantum mechanics, and quantum computation. Chapter 3 takes a deeper look into HQCAs and their working. Chapter 4 demonstrates the practical execution of HQCAs using Quantum Inspire’s quantum computing platform and SURF’s HPC center, analyzes the efficiency of different hybrid quantum-classical workflows, and

recommends changes to improve the efficiency. Finally, Chapter 5 concludes this work by reviewing the results and discussing possible future work.

2

Background

This chapter is focused on making the reader familiar with concepts used throughout this report. First, an introduction to computational complexity is given to establish a mathematical framework to describe the efficiency of computer algorithms. Second, the basic ideas of quantum mechanics are presented. Finally, an overview of the quantum circuit computational model is given along with a description of two important quantum algorithms to demonstrate the computing power of quantum computers.

2.1 Computational Complexity

There seems to be a fundamental limit to what problems computers can solve. Some problems seem to be inherently uncomputable: there exists no general solution that does not go into an infinite loop for certain inputs [15, 16]. On the other hand, there are limitations to the efficiency in which problems can be solved. This section will look at how one can define the computational efficiency of certain algorithms. That is, how many resources are required to solve a certain problem?

2.1.1 Big-O Notation

The time and space required by an algorithm generally grows as the size of the input grows. Because of this, it is traditional to describe the efficiency of an algorithm as a function of the size of its input [17]. This function describes the number of primitive operations it performs for a given input size. The notion of input size here depends on the context of the problem. For example, when computing the discrete Fourier transform, the input size likely refers to the dimension of the input vector. When talking about a problem like integer multiplication, however, it is more fitting to talk about the input size as the number of bits needed to represent the input in binary.

When analyzing the efficiency of algorithms, we look at the asymptotic growth for a given input size. Consider an algorithm that given input size n takes n^2 primitive operations to run and another algorithm that takes $500n^2 + \log n$ primitive operations to run. In big-O notation, both these algorithms are said to run in $O(n^2)$ time. That is, the number of primitive operations it performs scales quadratically with the input size. Constant factors are ignored as they become negligible as $n \rightarrow \infty$. While they are practically significant — an algorithm that runs in $O(n/2)$ performs half as many

primitive operations as an algorithm that runs in $O(n)$ — they are not relevant to asymptotic analysis, as both functions grow linearly with n .

Formally, if we have functions $f(n)$ and $g(n)$ such that $f(n)$ eventually grows slower than $Mg(n)$ for some real M as $n \rightarrow \infty$, we say

$$f(n) = O(g(n)). \quad (2.1)$$

For example, given $f(n) = 10n^3 + 7n^2$ and $g(n) = n^3$, $f(n)$ begins to grow slower than $Mg(n)$ with $M = 11$ for $n > 7$. Thus, g bounds f from above, and $f(n) = O(g(n)) = O(n^3)$. Some common big-O run times are shown in Table 2.1, along with their written name and an example. Throughout this report, algorithms that are bounded above by a polynomial (i.e. all run times until polynomial in Table 2.1) will be referred to as polynomial-time algorithms, and algorithms that are not bounded above by a polynomial will be referred to as superpolynomial-time algorithms.

Notation	Name	Example
$O(1)$	Constant	Accessing single element from array
$O(\log n)$	Logarithmic	Binary search
$O(n)$	Linear	Unstructured database search
$O(n \log n)$	Linearithmic	Fast Fourier Transform
$O(n^2)$	Quadratic	Insertion sort
$O(n^k)$	Polynomial	Gaussian elimination
$O(k^n)$	Exponential	Graph coloring
$O(n!)$	Factorial	Brute-force search traveling salesman problem

Table 2.1: Common big-O run times from fast to slow.

2.1.2 Turing Machines

The previous section described the measurement of computational efficiency as the number of primitive operations it performs for a given input size. This abstract definition can be extended by choosing a computational model in order to define what a primitive operation means. The standard computational model used for this is the Turing machine. It is chosen as computational model for the analysis of computational efficiency because of its simplicity and because it is able to simulate most physically realizable computational models with little overhead [18].

A Turing machine is an abstract machine that manipulates symbols from a work alphabet on a finite amount of one-way infinite length tapes divided into cells [16] (Figure 2.1). Along these tapes runs a tape head that can read and write one symbol at a time. The machine has a finite set of states, which the machine executes one at a time by loading them into the state register. At any time, the machine can be in one of the finite states. A state can be thought of as a rule with the following form:

$$(q_i, a) \mapsto (q_j, b, H), \quad (2.2)$$

where q_i and q_j are states, a and b are symbols from the work alphabet, and $H \in \{L, S, R\}$ decides how to move the tape head: one cell to the left (L), stay in the same position (S), or one cell to the right (R). These states as described in Equation 2.2 can be read as “in state q_i , if the read symbol is a , go to state q_j , write symbol b , and move the tape head to H .”

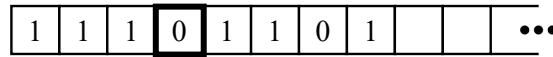


Figure 2.1: The tape of a single-tape Turing machine in an arbitrary state. Note that any multi-tape Turing machines can be efficiently simulated by a single-tape Turing machine [19], so complexity classes are not affected by changing between single-tape and multi-tape machines.

Everything that can be computed on models of computations we use these days can be computed on a Turing machine [20]. This hypothesis is known as the Church-Turing thesis. Related to the Church-Turing thesis is the extended Church-Turing thesis, which states that any physically realizable model of computation can be efficiently simulated on a Turing machine. That is, can a Turing machine simulate any model of computation in polynomial time? The quantum computational model brings doubt to this claim. It is known that quantum computers can efficiently simulate a Turing machine [21], but there appears to be no efficient algorithm for simulating a quantum computer on a Turing machine [22]. Furthermore, Arute, Arya, Babbush, *et al.* [23] experimentally demonstrated a quantum computer sampling from a probability distribution intractable by a classical computer. This gives us good reasons to believe that classical computers cannot efficiently simulate quantum computers, and that quantum computers are more computationally powerful for some problems.

2.1.3 Complexity Classes

Complexity classes are sets of computational problems that share some common feature with regard to the computational resources they need to solve some problem [18]. They are defined in terms of a type of computational problem, computational model, and a bounded resource such as time or space. In general, most complexity classes describe decision problems solvable by deterministic Turing machines — though many complexity classes are defined in terms of other types of problems and computational models. This report mainly focuses on complexity classes involving Turing machines and quantum computers.

The class P contains all decision problems solvable by a deterministic Turing machine in polynomial time. Problems that fall under this class are often referred to as tractable or easy problems [17]. The class NP (non-deterministic polynomial) contains all problems *verifiable* by a deterministic Turing machine in polynomial time. It is clear that $P \subseteq NP$, because if you can solve a problem in polynomial time, you can also verify it in polynomial time — you can verify it by solving it. A still unsolved and important question in computer science is whether $P = NP$. That is, can all problems that can be verified in polynomial time also be solved in polynomial time? A problem is said to be NP-hard when every problem in NP can be reduced into it in polynomial time, even if it is outside of NP. In term, if a problem is both NP-hard and in NP, the problem is NP-complete. NP-complete problems are considered to be the hardest problems in NP. By definition of NP-hardness, if a polynomial-time algorithm can be found for any NP-complete algorithm, then all problems in NP can be solved in polynomial time and $P = NP$.

In computer science, it is sometimes possible to speed up computation using randomness. These kinds of algorithms are referred to as probabilistic algorithms and are defined in terms of a probabilistic Turing machine. A probabilistic Turing machine is a Turing machine that can choose from a set of possible transitions according to some probability distribution when transitioning. The probabilistic equivalent of P is BPP (bounded-

error probabilistic polynomial time) and contains all decision problems solvable by a probabilistic Turing machine in polynomial time where a bounded error rate of $1/3$ is allowed. Since a probabilistic Turing machine can efficiently simulate a deterministic Turing machine, $P \subseteq BPP$. There are problems to be known in BPP and not in P, but the number of such problems is decreasing, and over the last decades a lot of evidence has been found towards $P = BPP$ [24–26].

How do quantum computers relate to these complexity classes? Quantum computers are probabilistic computational devices, and its complexity class equivalent to P can be defined by replacing the probabilistic Turing machine from BPP with a quantum computer.¹ The class BQP (bounded-error quantum polynomial time) consists of all decision problems solvable by a quantum computer in polynomial time where a bounded error rate of $1/3$ is allowed. It is known that there are NP problems that can be efficiently solved on a quantum computer like integer factorization, discrete logarithms, and quantum many-body simulation. As mentioned in the previous section, quantum computers can also solve all problems in P efficiently, so $P \subseteq BQP$. Furthermore, quantum computers are more powerful than classical probabilistic computers [27], giving $BPP \subseteq BQP$. How BQP relates to P and NP exactly is still unknown, however, it seems unlikely that $BQP = NP$ [28]. It also seems very unlikely that quantum computers will be able to solve NP-complete problems in polynomial time, or $BQP \cap NP\text{-complete} = \emptyset$.

The quantum analog of NP goes by the name of QMA (Quantum Merlin Arthur), and is related to BQP in the same way that NP is related to P. That is, QMA is the set of all decision problems verifiable by a quantum computer in polynomial time where a bounded error rate of $1/3$ is allowed. Similarly then, a problem is said to be QMA-hard when every problem in QMA can be reduced to it, and a problem is said to be QMA-complete if it is both QMA-hard and in QMA. Just like P and NP, it is clear that $BQP \subseteq QMA$, and $P \subseteq NP \subseteq QMA$.

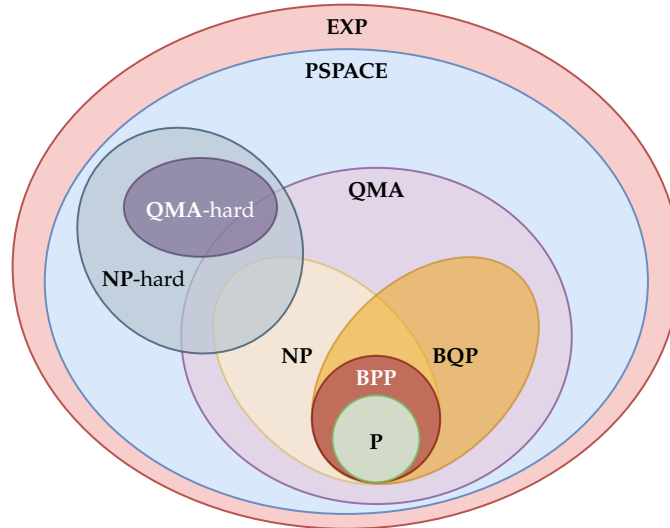


Figure 2.2: An overview of the hierarchy of some of the complexity classes discussed. This graphic assumes $P \neq NP$, $P \neq BPP$, and $P \subseteq BPP \subseteq BQP$. PSPACE is the space equivalent of P containing all problems that can be solved in polynomial space by a deterministic Turing machine, and EXP contains all problems that can be solved in exponential time by a deterministic Turing machine.

¹Note that quantum computers are not simply probabilistic Turing machines as will be shown in the following sections.

2.2 Quantum Mechanics

Quantum mechanics is the most complete description of the physical properties of nature on the atomic scale. It is also at the core of quantum computation and information. This section describes the necessary background of quantum mechanics required for understanding quantum computing.

2.2.1 Qubits

In classic information theory, the smallest unit of information is the bit. Quantum information is built upon an analogous concept: the quantum bit, or qubit. Qubits are physical objects that appear in nature on the scale of atoms and subatomic particles. A qubit can be any two-state quantum-mechanical system such as the spin of an electron, which can be spin up or down, or the polarization of a photon, which can be horizontally or vertically polarized. In this report, qubits will be treated as abstract mathematical objects as the physical realization of qubits is beyond the scope of this work. The state of a qubit is denoted as follows:

$$|\psi\rangle = \alpha_0|0\rangle + \alpha_1|1\rangle. \quad (2.3)$$

Quantum states are often described using Dirac notation $|\cdot\rangle$, which describes a column vector in \mathbb{C}^{2^n} . The states $\{|0\rangle, |1\rangle\}$ are the computational basis states which are defined as $(1 \ 0)^T$ and $(0 \ 1)^T$ respectively, and form an orthonormal basis for this vector space. The values $\alpha, \beta \in \mathbb{C}$ are the state's probability amplitudes, and cannot be examined directly.² When measuring a qubit, it collapses probabilistically to one of the basis states. The probability of measuring 0 is given by the absolute square $|\alpha_0|^2$, and the probability of measuring 1 is given by $|\alpha_1|^2$. As these values are probabilities, they should be normalized: $|\alpha_0|^2 + |\alpha_1|^2 = 1$. Formally, a qubit can be thought of as a unit vector in a two-dimensional Hilbert space.

A qubit differs from a classical bit in that it can be in a linear combination, or superposition of states. While a bit can be only be in the state 0 or 1, a qubit can be in one of infinitely many superpositions of states. However, the laws of quantum mechanics restrict direct access to the probability amplitudes of a state. Instead, when measuring a qubit, it collapses to basis state $|j\rangle$ with probability $|\alpha_j|^2$. For example, consider the state

$$|+\rangle = \frac{1}{\sqrt{2}} (|0\rangle + |1\rangle). \quad (2.4)$$

This state has equal probability of measuring 0 and 1, as $|1/\sqrt{2}|^2 = 1/2$. Note that measurement changes the state of a qubit: if the state from Equation 2.4 is measured as 1, the superposition is lost and the state becomes $|1\rangle$.

A helpful geometric interpretation of a qubit's state can be obtained by rewriting Equation 2.3 as

$$|\psi\rangle = e^{i\delta} \left(\cos \frac{\theta}{2} |0\rangle + e^{i\varphi} \sin \frac{\theta}{2} |1\rangle \right), \quad (2.5)$$

where $\delta, \theta, \varphi \in \mathbb{R}$. The global phase $e^{i\delta}$ can often be ignored, as $\forall \delta \in \mathbb{R} : |e^{i\delta}| = 1$, so it does not impact measurement outcome. Simplifying then, the state of a qubit can be written as

$$|\psi\rangle = \cos \frac{\theta}{2} |0\rangle + e^{i\varphi} \sin \frac{\theta}{2} |1\rangle. \quad (2.6)$$

²The field of quantum tomography focuses on recovering these values through multiple measurements, but this requires prior knowledge about the system [29].

Here, θ and φ define a point on the surface of a three-dimensional sphere referred to as the Bloch sphere (Figure 2.3). While this visualization is limited to a single qubit, it can

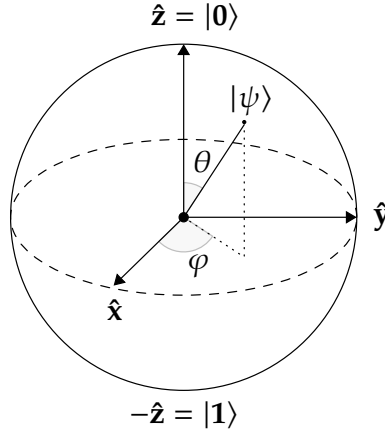


Figure 2.3: Bloch sphere representation of a qubit's state.

be a useful visual to build intuition. For example, the $|+\rangle$ state described in Equation 2.4 can be thought of as being exactly between $|0\rangle$ and $|1\rangle$ on the Bloch sphere.

The amount of probability amplitudes grows exponentially with the number of qubits: a n -qubit state has $N = 2^n$ amplitudes. Consider a two-qubit system which lives in a $2^2 = 4$ -dimensional Hilbert space spanned by the computational basis states $\{|00\rangle, |01\rangle, |10\rangle, |11\rangle\}$. This state is defined by the linear combination

$$|\psi\rangle = \alpha_0|00\rangle + \alpha_1|01\rangle + \alpha_2|10\rangle + \alpha_3|11\rangle. \quad (2.7)$$

Again, unlike classical bits who can only be in one state at a time, this state can be in a superposition of all four states. The normalization condition still applies for Equation 2.7: $\sum_{j=0}^3 |\alpha_j|^2 = 1$. Single-qubit states can be combined to form multi-qubit states by taking the tensor product of the two states. Given states $|\psi\rangle = \alpha_0|0\rangle + \alpha_1|1\rangle$ and $|\varphi\rangle = \beta_0|0\rangle + \beta_1|1\rangle$:

$$\begin{aligned} |\psi\rangle \otimes |\varphi\rangle &= \begin{pmatrix} \alpha_0 \\ \alpha_1 \end{pmatrix} \otimes \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \\ &= \begin{pmatrix} \alpha_0 \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \\ \alpha_1 \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \end{pmatrix} \\ &= \begin{pmatrix} \alpha_0\beta_0 \\ \alpha_0\beta_1 \\ \alpha_1\beta_0 \\ \alpha_1\beta_1 \end{pmatrix} \\ &= \alpha_0\beta_0|00\rangle + \alpha_0\beta_1|01\rangle + \alpha_1\beta_0|10\rangle + \alpha_1\beta_1|11\rangle. \end{aligned} \quad (2.8)$$

Often when describing multi-qubit states, the tensor product will be implied. The notations

$$|0\rangle \otimes |0\rangle = |0\rangle|0\rangle = |00\rangle \quad (2.9)$$

all describe state. The relative phases of α_0, α_1 and β_0, β_1 in Equation 2.8 are responsible for the quantum mechanical property of interference. When the phase of α_j and β_k is

the same, they will interfere constructively and increase the probability amplitude for that state. On the other hand, if α_j and β_k have opposite phases, they will interfere destructively and decrease the probability amplitude for that state.

Not all multi-qubit systems can be expressed as a tensor product of individual states as shown in Equation 2.8. Consider the following state:

$$|\Phi^+\rangle = \frac{1}{\sqrt{2}} (|00\rangle + |11\rangle). \quad (2.10)$$

This state cannot be expressed as a tensor product of two individual states, as that would imply $(\alpha_0\beta_0 = \alpha_1\beta_1 = 1/\sqrt{2}) \wedge (\alpha_0\beta_1 = \alpha_1\beta_0 = 0)$, which is a contradiction. States like $|\Phi^+\rangle$ are referred to as entangled states. Entanglement is the quantum phenomena of correlation in measurement outcomes. For example, when measuring the state $|\Phi^+\rangle$ from Equation 2.10, the only two possible measurement outcomes are 00 and 11. So by measuring one qubit, one also knows the state of the other qubit.

2.2.2 State Evolution

The evolution of a closed quantum system is described by a unitary transformation [30, Section 2.2].³ A state $|\psi\rangle$ at time t_1 is related to state $|\psi'\rangle$ at time t_2 by a unitary operator U :

$$|\psi'\rangle = U|\psi\rangle. \quad (2.11)$$

The unitary nature of these operators implies $UU^\dagger = U^\dagger U = I$, where † is the conjugate transpose and I the identity matrix. Single-qubit operators can be represented as 2×2 complex-valued unitary matrices. A common single-qubit operator is the Pauli-X operator which transforms a state $\alpha_0|0\rangle + \alpha_1|1\rangle$ to $\alpha_1|0\rangle + \alpha_0|1\rangle$. It is part of the set of Pauli matrices:

$$I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad Y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, \quad Z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}. \quad (2.12)$$

These matrices are ubiquitous in the study of quantum computation and information. Another useful and common operator is the Hadamard operator:

$$H = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, \quad (2.13)$$

which maps the computational basis states to an equal superposition state.

The Pauli X, Y and Z operators give rise to the rotation operators about the \hat{x} , \hat{y} , and \hat{z} axes when exponentiated:

$$R_x(\theta) = e^{-i\frac{\theta}{2}X} = \begin{pmatrix} \cos \frac{\theta}{2} & -i \sin \frac{\theta}{2} \\ -i \sin \frac{\theta}{2} & \cos \frac{\theta}{2} \end{pmatrix}, \quad (2.14)$$

$$R_y(\theta) = e^{-i\frac{\theta}{2}Y} = \begin{pmatrix} \cos \frac{\theta}{2} & -\sin \frac{\theta}{2} \\ \sin \frac{\theta}{2} & \cos \frac{\theta}{2} \end{pmatrix}, \quad (2.15)$$

$$R_z(\theta) = e^{-i\frac{\theta}{2}Z} = \begin{pmatrix} e^{-i\frac{\theta}{2}} & 0 \\ 0 & e^{i\frac{\theta}{2}} \end{pmatrix}. \quad (2.16)$$

³In this report, a perfectly closed system is assumed, even though in reality all systems interact somewhat with other systems. This is not an issue for practical quantum computers though due to quantum error correction, which can protect quantum information against noise from the external environment.

These operators can be thought of as rotating a qubit's state among its relative axis by an angle θ . Note that any single-qubit operator U can be decomposed into these operators, for example $U = R_z(\gamma)R_y(\beta)R_z(\alpha)$ where $\alpha, \beta, \gamma \in \mathbb{R}$.

Multi-qubit operators act on two or more qubits and are required for creating entangled states. Two common two-qubit operators are the CNOT and CZ operators. The CNOT operator can be thought of as a controlled-X operator, which applies a Pauli-X operation on the target qubit if the control qubit is in the $|1\rangle$ state. Equivalently, the CZ operator applies a Pauli-Z operation on the target qubit if the control qubit is $|1\rangle$. Their matrix representations are as follows:

$$\text{CNOT} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \quad \text{CZ} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix}. \quad (2.17)$$

Generally, a controlled- U operator applies U on the target qubit if the control qubit is $|1\rangle$.

2.2.3 Measurement

While the evolution of quantum states in a closed system is unitary, at some point the quantum state has to interact with the outside world. This is what measurement means: any interaction from an outside system with the quantum system. Formally, measurement is defined by a set of measurement operators $\{M_m\}$, where the probability of measuring the state m is given by

$$p(m) = \langle \psi | M_m^\dagger M_m | \psi \rangle. \quad (2.18)$$

With $\langle \psi | = |\psi\rangle^\dagger$, this equation can then be read as the inner product between $M_m|\psi\rangle$ and itself. This is just a generalization of the definition of measurement given in Section 2.2.1. For example, what is the probability of measuring 0 for an arbitrary state $|\psi\rangle = \alpha_0|0\rangle + \alpha_1|1\rangle$? Using $M_0 = |0\rangle\langle 0|$:

$$\begin{aligned} p(0) &= \langle \psi | M_0^\dagger M_0 | \psi \rangle \\ &= \langle \psi | 0 \rangle \langle 0 | 0 \rangle \langle 0 | \psi \rangle \\ &= \langle \psi | 0 \rangle \langle 0 | \psi \rangle \\ &= \begin{pmatrix} \alpha_0^* & \alpha_1^* \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \end{pmatrix} \\ &= \alpha_0^* \alpha_0 \\ &= |\alpha_0|^2. \end{aligned} \quad (2.19)$$

Similarly, using $M_1 = |1\rangle\langle 1|$ gives $p(1) = |\alpha_1|^2$. Measuring with the measurement operators $\{|0\rangle\langle 0|, |1\rangle\langle 1|\}$ is referred to as measuring in the computational basis.

After measurement, the state of the system can be described as follows:

$$|\psi\rangle = \frac{M_m|\psi\rangle}{\sqrt{p(m)}}. \quad (2.20)$$

Following the example from Equation 2.19 with $M_0 = |0\rangle\langle 0|$ and $p(0) = |\alpha_0|^2$, after

measuring 0 the system is in the state

$$\begin{aligned}
 |\psi\rangle &= \frac{|0\rangle\langle 0|\psi\rangle}{\sqrt{|\alpha_0|^2}} \\
 &= \frac{|0\rangle\alpha_0}{|\alpha_0|} \\
 &= \frac{\alpha_0}{|\alpha_0|}|0\rangle.
 \end{aligned} \tag{2.21}$$

Note that a global phase $e^{i\delta}$ shows up as the factor $\alpha_0/|\alpha_0|$. As mentioned in Section 2.2.1, the states $e^{i\delta}|0\rangle$ and $|0\rangle$ are considered equal up to the global phase factor.

A special class of measurements called projective measurements can sometimes be used to simplify calculations. A projective measurement is defined by an observable O , which is a Hermitian operator on the state space of the system. This observable has a spectral decomposition

$$O = \sum_j \lambda_j |\lambda_j\rangle\langle\lambda_j|, \tag{2.22}$$

where λ_j are eigenvalues of O and $|\lambda_j\rangle$ are the respective eigenstates. For example, the Pauli-Z operator can be thought of as an observable with the spectral decomposition

$$Z = 1|0\rangle\langle 0| - 1|1\rangle\langle 1|, \tag{2.23}$$

which has eigenstates $\{|0\rangle, |1\rangle\}$ with respective eigenvalues $\{1, -1\}$. As the observable Z has the computational basis states as eigenstates, a measurement of Z can also be thought of as measuring in the computational basis. With this definition of projective measurements, the expectation value of an observable O for state $|\psi\rangle$ can be defined as

$$\begin{aligned}
 \langle O \rangle_\psi &= \langle\psi|O|\psi\rangle \\
 &= \langle\psi|\left(\sum_j \lambda_j |\lambda_j\rangle\langle\lambda_j|\right)|\psi\rangle \\
 &= \sum_j \lambda_j \langle\psi|\lambda_j\rangle\langle\lambda_j|\psi\rangle \\
 &= \sum_j \lambda_j |\langle\lambda_j|\psi\rangle|^2.
 \end{aligned} \tag{2.24}$$

The expectation value is the sum of all possible outcomes (eigenvalues of O) weighted by their probability. Calculating the expectation value experimentally means preparing and measuring the state multiple times and is a useful way of extracting a deterministic quantity from a non-deterministic model of computation.

2.3 Quantum Computation

This section combines the fields of computer science and quantum mechanics to introduce the fundamental model of quantum computation: the quantum circuit model.

2.3.1 Quantum Circuits

Analogous to the classical circuit model, the quantum circuit model uses gates which act on data. In the classical circuit model boolean functions (logic gates) act on bits, while in the quantum circuit model quantum gates act on qubits. Most of the common quantum gates used in quantum computing were already described as operators in Section 2.2.2. To be more in line with the nomenclature of classical computing, from here on out these operators will be referred to as quantum gates in the context of quantum circuits.

In a quantum circuit time moves from left to right, where qubits are represented by wires on which gates can act, and classical bits are represented by double-lined wires. Usually the qubits are assumed to be instantiated to $|0\rangle$, unless noted otherwise. A list of frequently used quantum gates is shown in Table 2.2. Quantum gates are unitary and thus reversible, making the quantum circuit model a reversible model of computation. The inverse of a gate U is denoted as the conjugate transpose U^\dagger . By the definition of unitary $UU^\dagger = U^\dagger U = I$, applying the inverse U^\dagger after U essentially uncomputes U and vice versa. Gates whose conjugate transpose are equal to themselves are referred to as Hermitian. For example, H is Hermitian as $H^\dagger = H$ and thus $H^2 = I$.


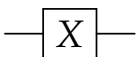
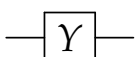

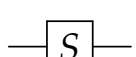
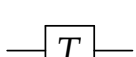
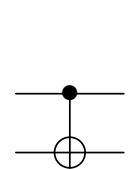
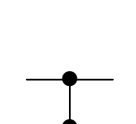
Gate name	Circuit symbol	Matrix representation
Hadamard		$\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$
Pauli-X		$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$
Pauli-Y		$\begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}$
Pauli-Z		$\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$
Phase (S)		$\begin{pmatrix} 1 & 0 \\ 0 & i \end{pmatrix}$
T		$\begin{pmatrix} 1 & 0 \\ 0 & e^{i\pi/4} \end{pmatrix}$
CNOT		$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$
CZ		$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix}$

Table 2.2: Frequently used quantum gates with their circuit symbol and matrix representation.

Figure 2.4 demonstrates a simple quantum circuit which creates the maximally entangled state $|\Phi^+\rangle = (|00\rangle + |11\rangle) / \sqrt{2}$ from Equation 2.10 and measures both qubits. This

state is one of four maximally entangled two-qubit states referred to as a Bell state. This

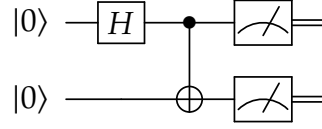


Figure 2.4: Quantum circuit for creating and measuring a Bell state $|\Phi^+\rangle$.

circuit does the following. The system starts in state $|\psi\rangle = |00\rangle$. A Hadamard gate is applied on the first qubit:

$$|\psi\rangle = \frac{1}{\sqrt{2}} (|0\rangle + |1\rangle) |0\rangle \quad (2.25)$$

$$= \frac{1}{\sqrt{2}} (|00\rangle + |10\rangle). \quad (2.26)$$

Then, a CNOT is applied with the first qubit as control and the second qubit as target, giving the final state

$$|\psi\rangle = \frac{1}{\sqrt{2}} (|00\rangle + |11\rangle). \quad (2.27)$$

This state is then measured, which will measure 00 or 11 with equal probability. Quantum circuits are often executed multiple times (multiple “shots”) to get a probability distribution of all possible outcomes.

Quantum gates can be applied in parallel on different qubits. Because of this, it is often more useful to talk about the depth of a quantum circuit rather than the number of gates. The depth of a quantum circuit is number of steps required for the quantum operations making up the circuit to run. For example, the quantum circuit in Figure 2.5 consists of 8 quantum gates, but has a depth of 4 (measurement gates are omitted from the depth calculation). The gates in every column are considered to be one time step as they can be executed in parallel. Instead of applying the first column of gates sequentially as $(I \otimes I \otimes S)(I \otimes X \otimes I)(H \otimes I \otimes I)|000\rangle$, it can be thought of as applying the single operator $(H \otimes X \otimes S)|000\rangle$.

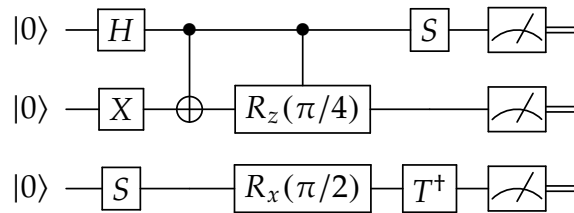


Figure 2.5: A quantum circuit consisting of 8 gates with a depth of 4.

In classical computations, copying bits is a common operation. However, in quantum computing the data one can copy is much more restricted. If a qubit is in a computational basis state, that state can be copied by a CNOT gate (Figure 2.6).⁴ However, trying to copy an arbitrary state $|\psi\rangle = \alpha_0|0\rangle + \alpha_1|1\rangle$ using the circuit in Figure 2.6 gives

$$|\psi\rangle = \alpha_0|00\rangle + \alpha_1|11\rangle. \quad (2.28)$$

⁴More formally, if a quantum state is in a basis state of a known orthogonal basis, the state can be copied.

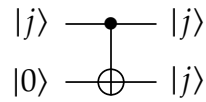


Figure 2.6: Copying a computational basis state using CNOT where $j \in \{0, 1\}$.

This state does not contain two copies of $|\psi\rangle$ (unless $\alpha_0\alpha_1 = 0$ as is true with computational basis states), which should have the form

$$|\psi\rangle|\psi\rangle = \alpha_0^2|00\rangle + \alpha_0\alpha_1|01\rangle + \alpha_1\alpha_0|10\rangle + \alpha_1^2|11\rangle. \quad (2.29)$$

In fact, is it impossible to copy an unknown quantum state. There exists no solution for a unitary operator U that does the transformation

$$\alpha_0|00\rangle + \alpha_1|01\rangle \xrightarrow{U} \alpha_0^2|00\rangle + \alpha_0\alpha_1|01\rangle + \alpha_1\alpha_0|10\rangle + \alpha_1^2|11\rangle. \quad (2.30)$$

This property is known as the no-cloning theorem, and is a fundamental limit of quantum information. Note that this holds for unknown quantum states. If one has the circuit to prepare $|\psi\rangle$, you could simply create $|\psi\rangle|\psi\rangle$ by executing the circuit on a second qubit.

2.3.2 Universal Gate Sets

A universal gate set is defined as a finite set of gates that can be used to represent any other gate. In the classical circuit model, NAND is a universal gate with which all other gates can be represented. Equivalently, universal quantum gate sets are sets of quantum gates to which any quantum operation can be reduced to. A common universal gate set is $\{\text{CNOT}, H, S, T\}$. Any quantum operation can be reduced to a combination of CNOT, H , S , and T gates. For example, the three-qubit equivalent of the CNOT gate, the Toffoli gate, can be decomposed into these gates as shown in Figure 2.7. The Solovay-Kitaev theorem shows that any universal gate set can simulate any other universal gate set efficiently [31], so the choice of universal gate set does not impact the asymptotic efficiency of a quantum computer. It turns out that almost any set of one- and two-qubit gates is universal — in fact, Shi [32] shows that just the Toffoli and Hadamard gates are universal.

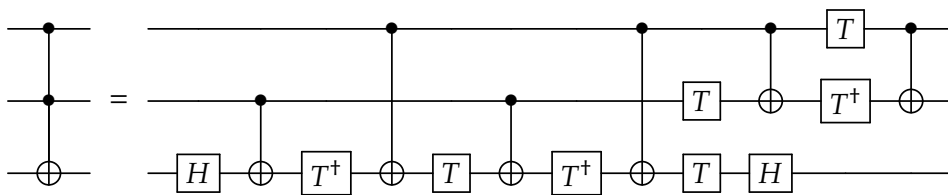


Figure 2.7: Decomposition of the Toffoli gate using CNOT, H , S and T gates.

2.3.3 Quantum Algorithms

To demonstrate the quantum circuit computational model's ability to solve certain problems more efficiently than classical computers, two important quantum algorithms which provide a speedup over their classical counterparts are now reviewed: the quantum Fourier transform (QFT) and Grover's algorithm.

Quantum Fourier Transform

The quantum Fourier transform (QFT) is the quantum analogue of the *inverse* discrete Fourier transform (DFT). The DFT is ubiquitous in the fields of engineering and computer science. In short, it transforms a signal in the time domain to the frequency domain, and vice versa with the inverse DFT. For $N = 2^n$ amplitudes, the QFT uses $O(n^2)$ gates, where n is the amount of qubits. The best-known classical algorithm for computing the DFT uses $O(n2^n)$ gates. That is, the QFT is exponentially more efficient than the best-known classical algorithm. However, the applications of the QFT are limited. The result of the QFT are stored in the amplitudes of the quantum state. These amplitudes cannot be directly accessed by measurement, so there is no way to extract the results. On the other hand, the QFT is the backbone of the phase estimation algorithm, which in turn can be used to solve problems including the order-finding problem and the factoring problem efficiently.

The inverse DFT is a linear transformation on an input vector of N complex values (x, \dots, x_N) . It transforms a element x_k as follows:

$$x'_k = \frac{1}{\sqrt{N}} \sum_{j=0}^{N-1} x_j e^{\frac{2\pi i}{N} jk}. \quad (2.31)$$

Equivalently, the QFT is defined as a unitary operator which acts on a computational basis state $|j\rangle$ as following:

$$|j\rangle \rightarrow \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} e^{\frac{2\pi i}{N} jk} |k\rangle. \quad (2.32)$$

To describe the quantum circuit for the QFT, the following definition is useful:

$$R_k = \begin{pmatrix} 1 & 0 \\ 0 & e^{\frac{2\pi i}{2^k}} \end{pmatrix}. \quad (2.33)$$

Some useful identities are $R_1 = Z$, $R_2 = S$, and $R_3 = T$. Using this gate, the general quantum circuit for the QFT is shown in Figure 2.8. The gates at the end of the circuit are swap gates, which swap two qubits. This is necessary to get the result in the correct order.

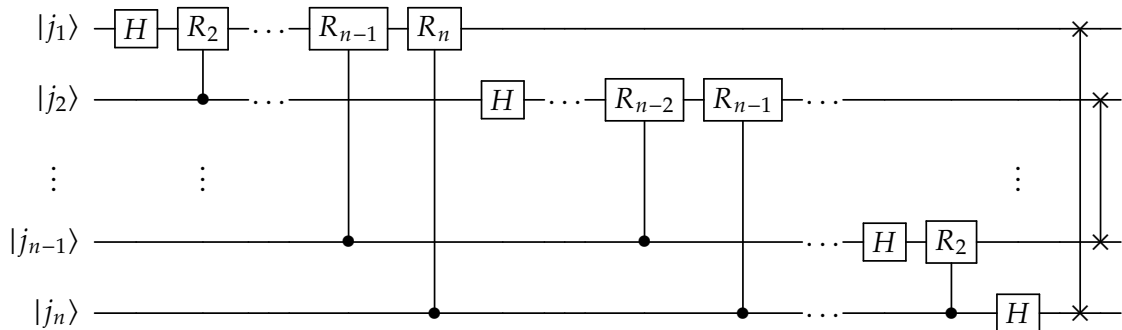


Figure 2.8: General circuit for the QFT on 2^n amplitudes with n qubits.

The QFT circuit works as follows. With $N = 2^n$ where n is the number of qubits, the computational basis exists of the states $\{|0\rangle, \dots, |N-1\rangle\}$. We use the notation $|j\rangle$ for a

decimal number j in the binary representation j_1, \dots, j_m , e.g. $|15\rangle = |j_1 j_2 j_3 j_4\rangle = |1111\rangle$. Furthermore, the following notation is used to represent fractional binary numbers:

$$[0.j_1 \dots j_m] = \sum_{k=1}^m j_k 2^{-k}. \quad (2.34)$$

For example, binary $[0.101]$ is $1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} = 5/8$. Using these definitions, the QFT can be represented as follows:

$$\text{QFT}|j_1 \dots j_n\rangle = \frac{1}{\sqrt{N}} \left(|0\rangle + e^{2\pi i [0.j_n]} |1\rangle \right) \left(|0\rangle + e^{2\pi i [0.j_{n-1} j_n]} |1\rangle \right) \dots \left(|0\rangle + e^{2\pi i [0.j_1 j_2 \dots j_n]} |1\rangle \right). \quad (2.35)$$

The transformed values end up in the relative phases of the amplitudes. Again, these relative phases cannot be extracted. However, the QFT is an essential part of important quantum algorithms: it is at the heart of the quantum phase estimation algorithm (QPEA), which can be used to factor integers [4], solve linear systems of equations [33], and more [30, Section 5.4].

Grover's Algorithm

Grover's algorithm is a quantum algorithm for unstructured database search. Suppose you want to search through a search space of n numbers. Given a function $f(x)$ that returns 1 for the number ω that you are looking for, and 0 for every other number:

$$f(x) = \begin{cases} 1 & \text{if } x = \omega \\ 0 & \text{if } x \neq \omega \end{cases} \quad (2.36)$$

The goal is to find a x so that $f(x) = 1$. On a classical computer, the best we can do is an exhaustive search of the entire search space, which takes $O(n)$ function calls. Using a quantum computer, this can be done with $O(\sqrt{n})$ function calls using Grover's algorithm [5]. Grover's algorithm is also asymptotically optimal: Zalka [34] showed that $O(\sqrt{n})$ is the best we can do for unstructured database search.

A common problem in computer science is examining a large number of different possibilities to see which of them satisfy a given condition. This is the same problem that Grover's algorithm looks to solve. Usually, there is some structure in the problem which allows for more efficient algorithms. For example, if a list of numbers is sorted, one can find a number in $O(\log n)$ using binary search. For hard problems that do not have such obvious structure, Grover's algorithm can be used to achieve a quadratic speedup.

The function $f(x)$ from Equation 2.36 is sometimes referred to as a black box function, as its internal workings is not of interest. Grover's algorithm is a general algorithm which can be applied to many algorithms that use search heuristics. The quantum equivalent of a black box function is called an oracle, and has the following form:

$$U_\omega |x\rangle = \begin{cases} -|x\rangle & \text{if } x = \omega \\ |x\rangle & \text{if } x \neq \omega \end{cases} \quad (2.37)$$

where $|x\rangle$ is any computational basis state, and ω is the bit string to find. That is, it marks the solution to the search problem by shifting the phase of the solution. An oracle from Equation 2.37 can be defined in terms of a black box function:

$$U_\omega |x\rangle = (-1)^{f(x)} |x\rangle, \quad (2.38)$$

which has the following matrix representation:

$$U_\omega = \begin{pmatrix} (-1)^{f(0)} & 0 & \cdots & 0 \\ 0 & (-1)^{f(1)} & \cdots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & 0 & \cdots & (-1)^{f(2^n-1)} \end{pmatrix}. \quad (2.39)$$

With the oracle defined, the quantum circuit of Grover's algorithm is shown in Figure 2.9. The quantum state can be thought of as having two registers: the target qubits $|q_1 \dots q_n\rangle$ and the oracle qubit $|y\rangle$. The oracle qubit $|y\rangle$ starts in the state $H|1\rangle = (|0\rangle - |1\rangle)/\sqrt{2}$. When the oracle is applied to $|q_1 \dots q_n\rangle(|0\rangle - |1\rangle)/\sqrt{2}$ and x is not a solution, the state stays the same. However, if the oracle is applied to $|q_1 \dots q_n\rangle(|0\rangle - |1\rangle)/\sqrt{2}$ and x is a solution, the state becomes $-|q_1 \dots q_n\rangle(|0\rangle - |1\rangle)/\sqrt{2}$. After the oracle is applied, the solution is marked with a negative phase. To extract this solution, a procedure called amplitude amplification is used. Amplitude amplification can be thought of as inversion about the mean.

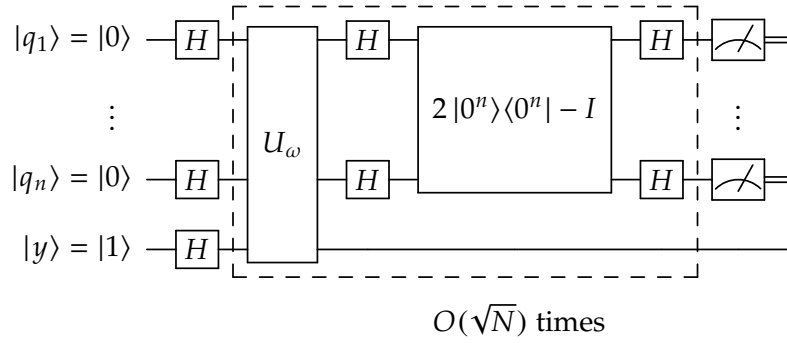
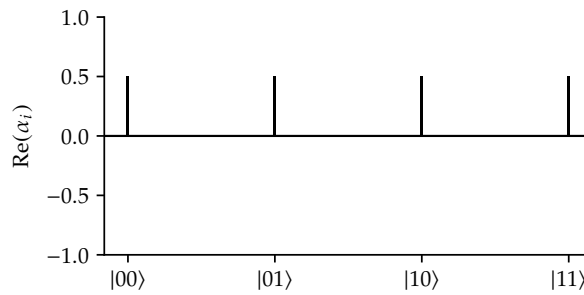
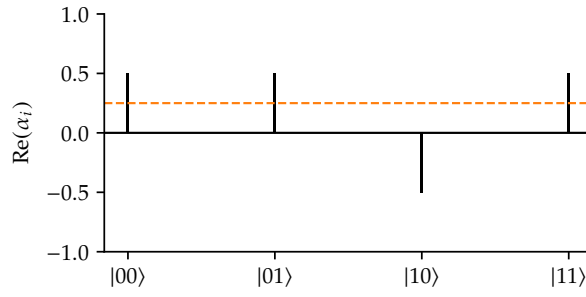


Figure 2.9: Quantum circuit representation of Grover's algorithm.

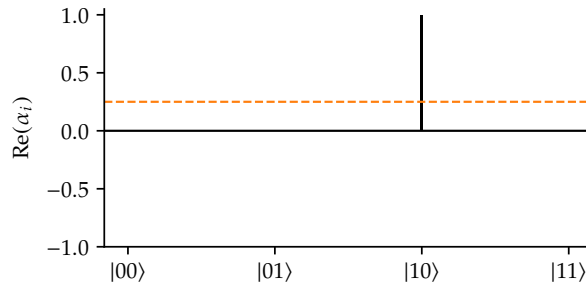
Consider a search space of $N = 4$ with two qubits where we are looking for $|\omega\rangle = |10\rangle$. The amplitudes of the search space start out in an equal superposition:



After applying the oracle U_ω the solution state $|10\rangle$ has its phase flipped:



Note that measuring this state will still give each state with equal probabilities. Using the diffusion operator $H^{\otimes n} (2 |0^n\rangle\langle 0^n| - I) H^{\otimes n}$, the amplitudes are inverted about their mean (displayed as the dotted orange line):



The application of the oracle and amplitude amplification is often referred to as a Grover iteration. For this example, one Grover iteration was enough to give a probability of 1 of measuring the correct solution. Generally, the upper bound on the number of required iterations R is as follows:

$$R \leq \left\lceil \frac{\pi}{4} \sqrt{\frac{M}{N}} \right\rceil, \quad (2.40)$$

where M is the number of solutions and $N = 2^n$ is the size of the search space.

3

Hybrid Quantum-Classical Algorithms

At the end of the previous chapter, two important quantum algorithms were reviewed: the QFT and Grover's algorithm. These quantum algorithms and the family of algorithms built on them require large-scale quantum computers with many qubits, low error rates, and high coherence times. The NISQ devices that are being built now and in the near future will not be able to run these algorithms due to their limited qubit count, limited coherence times, and high error rates. To this end, hybrid quantum-classical algorithms (HQCAs) that utilize both classical and quantum resources are being researched and developed. This chapter will first introduce the basic concepts of an important family of HQCAs followed by a review of well-known and promising HQCAs.

3.1 Variational Hybrid Quantum-Classical Algorithms

Hybrid quantum-classical algorithms take into account the limited number of qubits, limited connectivity of qubits, high error rates, and limited coherence times of NISQ devices. These algorithms often make use of the variational method which consists of preparing an initial trial state $|\psi(\theta)\rangle$ parameterized by real-valued parameters θ , and finding the parameters for which the expectation value of some observable is the lowest. This family of HQCAs is sometimes also referred to as variational quantum algorithms. The rest of this chapter will be focused on variational HQCAs, as this is the main focus of research regarding HQCAs.

The general structure of a variational HQCA is shown in Figure 3.1. While the specific implementation differs between algorithms, they all share some basic elements. First, a cost function C which encodes the solution to the problem needs to be defined. Second, an ansatz needs to be chosen. The ansatz defines the operations of the quantum circuit which are parameterized by the trainable parameters θ . The parameters θ are then trained in a quantum-classical loop to solve the optimization problem

$$\theta_{\text{opt}} = \arg \min_{\theta} C(\theta). \quad (3.1)$$

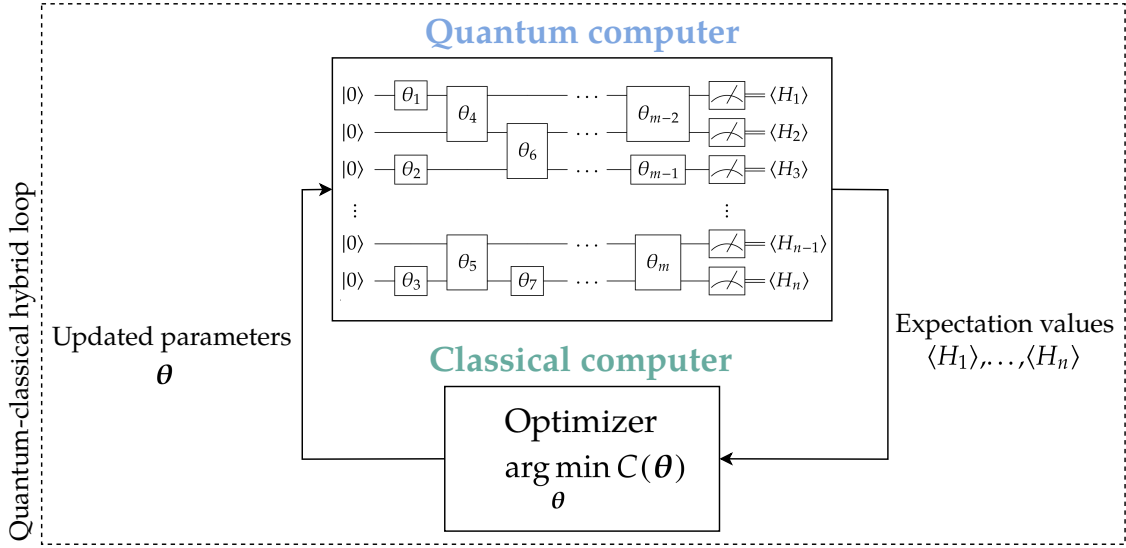


Figure 3.1: The general structure of a variational HQCA. A parameterized quantum state $|\psi(\theta)\rangle$ is prepared through a set of parameterized quantum gates $U_L(\theta_L) \dots U_2(\theta_2)U_1(\theta_1)$. These parameters are classically optimized according to a cost function $C(\theta)$. This quantum-classical loop is repeated until convergence.

3.1.1 Cost Function

In general, the cost function has the form

$$C(\theta) = \sum_k f_k \left(\langle 0^n | U(\theta)^\dagger H_k U(\theta) | 0^n \rangle \right), \quad (3.2)$$

where f_k is a function encoding the problem at hand, $U(\theta)$ is the quantum circuit parameterized by real values θ , and H_k is an observable. A key point of variational HQCAs is that a quantum computer is used to estimate $C(\theta)$ while using a classical computer to optimize the parameters θ . Furthermore, the cost function $C(\theta)$ should be efficient to estimate by performing measurements on a quantum computer. For the variational HQCA to have a quantum advantage, the cost function should also not be efficiently computable with a classical computer.

3.1.2 Ansätze

The ansatz defines the quantum circuit $U(\theta)$ that is used to prepare the quantum state $|\psi(\theta)\rangle$ parameterized by real-valued parameters θ . To keep the algorithm realizable on NISQ devices, it is important to keep the parameterized quantum circuit (the ansatz) relatively shallow in depth. The challenge then is to choose an ansatz that is shallow in depth but has high expressibility. The expressibility of a quantum circuit is its ability to generate states that are well representative of the Hilbert space [35]. The general form of an ansatz can be described as the product of sequentially applied unitaries:

$$U(\theta) = U_L(\theta_L) \dots U_2(\theta_2)U_1(\theta_1), \quad (3.3)$$

where

$$U_l(\theta_l) = \prod_m e^{-i \frac{\theta_m}{2} H_m} W_m. \quad (3.4)$$

Here H_m is a Hermitian operator and W_m is a non-parameterized unitary. Equation 3.4 can simply be thought of as alternating between parameterized and non-parameterized unitaries.

The specific structure of the ansatz usually depends on the problem, as one can sometimes employ information about the problem to tailor an ansatz [36]. Such ansätze are referred to as problem-inspired ansätze. On the other hand, there are problem-agnostic ansätze that are used when no information is available for tailoring a specific ansatz.

3.1.3 Optimization

With the cost function and ansatz defined, the next step in a variational HQCA is to train the parameters θ using a classical optimization algorithm to solve the optimization problem from Equation 3.1. The two main types of optimization methods to consider are gradient-based methods and gradient-free methods. The former methods involve calculating first- or higher-order derivatives of the cost function, while the latter methods use only function evaluations. In general, the computational cost of gradient-based methods is larger, but the rate of convergence is significantly improved.

While gradient-free methods are straightforward to use for optimizing the quantum circuit parameters, gradient-based methods are more involved. Consider a parameterized gate of the form $G(\theta) = e^{-i\theta/2H}$ from the general ansatz description in Equation 3.4. Using the parameter-shift rule as described by Schuld, Bergholm, Gogolin, *et al.* [37], the derivative of the gate G with respect to a parameter θ can be described as follows:

$$\frac{\partial G}{\partial \theta} = \frac{1}{2} \left[G\left(\theta + \frac{\pi}{2}\right) - G\left(\theta - \frac{\pi}{2}\right) \right]. \quad (3.5)$$

That is, the derivative of G with respect to a parameter θ can be computed analytically by evaluating the quantum circuit twice with shifted parameters $\theta \pm \pi/2$.

For example, the derivative of a cost function $C(\theta) = \sum_k \langle 0^n | U(\theta)^\dagger H_k U(\theta) | 0^n \rangle$ with respect to a parameter θ_l can be estimated as follows:

$$\begin{aligned} \frac{\partial C}{\partial \theta_l} = \sum_k \frac{1}{2} & \left[\langle 0^n | U\left(\theta + \frac{\pi}{2} e_l\right)^\dagger H_k U\left(\theta + \frac{\pi}{2} e_l\right) | 0^n \rangle \right. \\ & \left. - \langle 0^n | U\left(\theta - \frac{\pi}{2} e_l\right)^\dagger H_k U\left(\theta - \frac{\pi}{2} e_l\right) | 0^n \rangle \right], \end{aligned} \quad (3.6)$$

where e_l is a vector with 1 as its l -th element and 0 for all other elements. Note that calculating ∂C with m parameters θ requires you to evaluate the quantum circuit twice for every parameter, for a total of $2m$ quantum circuit executions per gradient evaluation. The parameter-shift rule can be further generalized to calculate higher-order derivatives of quantum gates as demonstrated by Mari, Bromley, and Killoran [38], allowing the use of optimization methods that use higher-order derivatives such as Newton's method.

A common way to perform gradient-based optimization is through gradient descent based algorithms. Gradient descent uses first-order derivatives for finding a local minimum of a differentiable function. Intuitively, it takes repeated steps in the opposite direction of the gradient of the function until the gradient is zero (Figure 3.2). Formally, for a set of parameters θ_t and a cost function $C(\theta)$, the updated parameters θ_{t+1} after one iteration of gradient descent are as follows:

$$\theta_{t+1} = \theta_t - \eta \nabla C(\theta_t), \quad (3.7)$$

where η is the step size which determines how big of a step the algorithm makes each iteration. Finding the optimal step size is a practical challenge: if it is too large the optimization algorithm will diverge, and if it is too small it will take a long time to converge. Methods have been developed to improve the gradient descent algorithm — for example, the gradient descent based optimization algorithm Adam adapts the step size during optimization to increase the efficiency and accuracy of solutions [39].

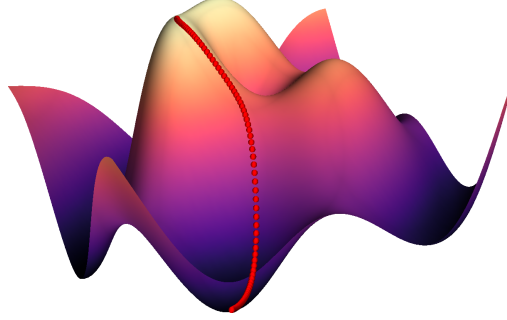


Figure 3.2: Visualization of gradient descent finding a local minimum of a two-dimensional function. The red dots represent the value of the cost function at a given iteration. The gradient descent algorithm gradually finds a local minimum by moving in the opposite direction of the gradient of the cost function.

3.2 Variational Quantum Eigensolver

The variational quantum eigensolver (VQE) was the first proposed variational HQCA. It aims to solve the problem of finding eigenvalues of a Hermitian matrix H . The problem of finding eigenvalues for large matrices is ubiquitous in quantum chemistry, but potential applications range from determining the results of search engines [40] to designing new materials and drugs [41]. A previously known quantum algorithm called the quantum phase estimation algorithm (QPEA) solves this problem exponentially faster than the best known classical algorithm [42]. However, the quantum resources required to run the QPEA greatly exceed the resources available in the near term. Peruzzo, McClean, Shadbolt, *et al.* [43] proposed the VQE as alternative which can run on NISQ devices.

Formally, the VQE estimates the minimum eigenvalue λ_{\min} associated with eigenstate $|\lambda_{\min}\rangle$ for a Hermitian operator H . The cost function is defined as the expectation value of H over a trial state $|\psi(\theta)\rangle = U(\theta)|\psi_0\rangle$ where $U(\theta)$ is the ansatz and $|\psi_0\rangle$ an initial guess state:

$$C(\theta) = \langle \psi(\theta) | H | \psi(\theta) \rangle. \quad (3.8)$$

Minimizing this cost function estimates the minimum eigenvalue of H , as the Rayleigh-Ritz variational principle provides the following bound [44]:

$$\lambda_{\min} \leq \langle \psi(\theta) | H | \psi(\theta) \rangle. \quad (3.9)$$

That is, a trial state $|\psi(\theta)\rangle$ will always give an expectation value larger than or equal to the minimum eigenvalue. In the context of quantum chemistry, where H describes the Hamiltonian of a system, finding the lowest eigenvalue of the Hamiltonian translates to finding the ground state and ground state energy of that system. By selecting an initial guess state, calculating the expectation value of H , and optimizing the cost function $C(\theta)$, one can then approximate the ground state energy of a Hamiltonian to arbitrary precision.

The problem of estimating the ground state energy is a specific instance of the k -local Hamiltonian problem, which looks to find the ground state energy of a k -local Hamiltonian. A k -local Hamiltonian is a Hamiltonian that can be expressed as a sum of Hamiltonians $H = \sum_j H_j$ where each H_j is a Hermitian operator acting non-trivially on at most k qubits [45]. This problem was shown to be QMA-complete for $k \geq 2$, while the 1-local Hamiltonian problem is in P [46]. However, quantum computers can still offer an advantage above classical computers for this problem. Any Hamiltonian H can be represented as a linear combination of Pauli operators:

$$H = \sum_j h_j P_j \quad (3.10)$$

$$= \sum_j h_j \prod_k P_j^{(k)}, \quad (3.11)$$

where $P_j^{(k)}$ is a Pauli operator $\{I, X, Y, Z\}$ which acts on qubit k , and h_j is a real constant. From the linearity of quantum observables, it follows that for some state $|\psi\rangle$

$$\langle\psi|H|\psi\rangle = \sum_j h_j \langle\psi|P_j|\psi\rangle \quad (3.12)$$

$$= \sum_j h_j \langle\psi| \prod_k P_j^{(k)} |\psi\rangle. \quad (3.13)$$

The problem of estimating $\langle\psi|H|\psi\rangle$ can be done efficiently on a quantum computer, while being a hard problem to solve classically [43, 47, 48].

In order to use the VQE to estimate ground state energies of molecules, the energy of the electrons and nuclei in a molecule need to be encoded as a Hamiltonian. The standard form of such Hamiltonian is as follows:

$$H = - \sum_i \frac{\nabla_{\mathbf{R}_i}^2}{2M_i} - \sum_i \frac{\nabla_{\mathbf{r}_i}^2}{2} - \sum_{i,j} \frac{Z_i}{|\mathbf{R}_i - \mathbf{r}_j|} + \sum_{i,j>i} \frac{Z_i Z_j}{|\mathbf{R}_i - \mathbf{R}_j|} + \sum_{i,j>i} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|}, \quad (3.14)$$

where M_i , Z_i , and \mathbf{R}_i denote the mass, atomic number, and position of nucleus i , and \mathbf{r}_i is the position of electron i [6]. Given the fact that the nuclei are much heavier than the electrons, we can apply the Born-Oppenheimer approximation [49] to separate the Hamiltonian into electronic and nuclear terms. We are mostly interested in the electronic structure of the molecule, so we can use the Born-Oppenheimer approximation to write the electronic Hamiltonian as

$$H = - \sum_i \frac{\nabla_{\mathbf{r}_i}^2}{2} - \sum_{i,j} \frac{Z_i}{|\mathbf{R}_i - \mathbf{r}_j|} + \sum_{i,j>i} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|}. \quad (3.15)$$

This form of Hamiltonian is often referred to as the first quantized formulation of quantum chemistry [50, Appendix A]. In order to represent electronic Hamiltonians on a quantum computer, it needs to be mapped from an operator acting on indistinguishable fermions to an operator acting on distinguishable qubits. There are several approaches to this [51], but the most widely used approach in quantum chemistry uses the second quantized formulation [6]. The second quantized Hamiltonian must then be mapped into qubits in order to be implemented on a quantum computer. The most common mappings for this are the Jordan-Wigner [52] and Bravyi-Kitaev [53] transformations.

The exact details on the different approaches are beyond the scope of this work, but a detailed overview can be found in [6].

Implementations of the VQE can often benefit from using a problem-inspired ansatz. A popular choice is the unitary coupled cluster (UCC) ansatz, which is a cousin of the coupled-cluster method used in quantum chemistry for describing many-body systems [54]. The UCC ansatz is parameterized by a polynomial number of parameters, while there is no known efficient classical implementation [55]. Furthermore, the UCC ansatz is believed to provide better accuracy than classical coupled cluster methods [56, 57]. Despite these advantages, the number of parameters required by the UCC ansatz might be too large to allow for practical calculations of large molecules [58]. Alternative ansätze have been proposed to reduce the quantum resources required for running the VQE on larger molecules [48, 59, 60].

3.3 Quantum Approximate Optimization Algorithm

The quantum approximate optimization algorithm (QAOA) is a variational HQCA that produces approximate solutions for combinatorial optimization problems [61]. Combinatorial optimization is focused on finding an optimal solution in a finite set of potential solutions. The number of potential solutions is often too large to scan one by one and select the best option, so more efficient methods are required. Most combinatorial optimization problems are known to be NP-complete [62]. The quantum approximate optimization algorithm (QAOA) depends on an integer $p \geq 1$, for which the quality of the approximation improves as p increases. Even for $p = 1$, the QAOA is not efficiently simulatable by classical computers [63]. The QAOA has received attention from the scientific community in the last few years because of its classical intractability, wide range of applications, and promising performance guarantees.

The QAOA looks to solve combinatorial optimization problems, which are specified by n bits and m clauses. A clause is a constraint on a subset of the bits which is satisfied for certain assignment of those bits, and not satisfied for other assignments. The goal is to find a n -bit string $z = \{0, 1\}^n$ that maximizes the classical objective function

$$L(z) = \sum_{\alpha=1}^m L_{\alpha}(z), \quad (3.16)$$

where $L_{\alpha}(z) = 1$ if z satisfies clause α and $L_{\alpha}(z) = 0$ if z does not satisfy clause α . The QAOA finds a bit string z' for which the approximation ratio $L(z')/L_{\max}$ is large, where $L_{\max} = \max_z L(z)$. The objective function L can be mapped to a Hamiltonian that is diagonal in the computational basis:

$$H_L = \sum_{\alpha=1}^m H_{L_{\alpha}}, \quad (3.17)$$

where

$$H_{L_{\alpha}} = \sum_{z \in \{0,1\}^n} L_{\alpha}(z) |z\rangle\langle z| \quad (3.18)$$

$$= \begin{pmatrix} L_{\alpha}(0, \dots, 0) & 0 & \cdots & 0 \\ 0 & L_{\alpha}(0, \dots, 1) & \cdots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & 0 & \cdots & L_{\alpha}(1, \dots, 1) \end{pmatrix}. \quad (3.19)$$

The operator H_L encodes the objective function and is referred to as the problem Hamiltonian.

The QAOA uses a problem-inspired ansatz which is defined in terms of two alternating unitaries, inspired by the quantum adiabatic algorithm [61]. First, we define a unitary parameterized by $\gamma \in [0, 2\pi)$ using the problem Hamiltonian:

$$U_L(\gamma) = e^{-i\gamma H_L} = \prod_{\alpha=1}^m e^{-i\gamma H_{L_\alpha}}. \quad (3.20)$$

Next, we define a mixer Hamiltonian H_B :

$$H_B = \sum_{j=0}^{n-1} X^{(j)}, \quad (3.21)$$

where $X^{(j)}$ is the Pauli-X gate acting on qubit j . Then, we define a mixer unitary parameterized by $\beta \in [0, 2\pi)$:

$$U_B(\beta) = e^{-i\beta H_B} = \prod_{j=0}^{n-1} e^{-i\beta X^{(j)}}. \quad (3.22)$$

Note that $e^{-i\beta X}$ is just a $R_x(2\beta)$ operation as defined in Section 2.2.2. Using these unitaries, the ansatz can be described as alternating applications of the problem and mixer Hamiltonian p times using different parameters:

$$U(\gamma, \beta) = \underbrace{U_B(\beta_p)U_L(\gamma_p) \dots U_B(\beta_1)U_L(\gamma_1)}_{p \text{ layers}}. \quad (3.23)$$

This unitary depends on $2p$ real parameters $\gamma = (\gamma_1, \dots, \gamma_p)$ and $\beta = (\beta_1, \dots, \beta_p)$. Note that not all quantum hardware can prepare this ansatz efficiently, so specific hardware-efficient ansätze are being proposed [48, 64].

Given these definitions, the QAOA works as follows. The system starts in the equal superposition state

$$|s\rangle = H^{\otimes n} |0\rangle^{\otimes n} = \frac{1}{\sqrt{N}} \sum_{j=0}^{N-1} |j\rangle. \quad (3.24)$$

A trial state $|\gamma, \beta\rangle$ is then obtained by applying the ansatz to $|s\rangle$:

$$|\gamma, \beta\rangle = U(\gamma, \beta)|s\rangle \quad (3.25)$$

$$= U_B(\beta_p)U_L(\gamma_p) \dots U_B(\beta_1)U_L(\gamma_1)|s\rangle, \quad (3.26)$$

and the cost function to optimize is as follows:

$$C_p(\gamma, \beta) = \langle \gamma, \beta | H_L | \gamma, \beta \rangle. \quad (3.27)$$

The subscript p is added to the cost function's definition to indicate the number of layers. Unlike the VQE which minimizes the cost function, we are interested in maximizing C_p to find the bit string which satisfies the most clauses. We look to solve the optimization problem

$$M_p = \max_{\gamma, \beta} C_p(\gamma, \beta). \quad (3.28)$$

By measuring $|\gamma, \beta\rangle$ multiple times in the computational basis we get a bit string z with high probability for which $L(z)$ is near or greater than M_p . The approximation ratio gets better as p increases, and as $p \rightarrow \infty$, $M_p = L_{\max}$.

4

Practical Hybrid Quantum-Classical Computing

This chapter is focused on the practical execution of HQCAs using the Quantum Inspire quantum computing platform and SURF's HPC center. It will start with describing the relevant platforms, followed by reviewing and analyzing how these platforms are used to execute HQCAs in different workflows. We will then explore and describe different methods to increase the efficiency of the execution of HQCAs.

4.1 Quantum Inspire

Quantum Inspire is a full-stack quantum computing platform that QuTech launched in 2020 to make quantum systems available to the general public for exploratory research [65]. Users can run quantum circuits on different back-ends through the Quantum Inspire web editor or by using the Python software development kit (SDK). The common quantum assembly language (cQASM) [66] is used for describing quantum circuits, but the popular quantum computing frameworks ProjectQ [67] and Qiskit [68] are also supported by the Python SDK. Quantum Inspire supports the execution of quantum circuits on real QPUs and through quantum simulation. An overview of the Quantum Inspire workflow and available device back-ends is shown in Figure 4.1.

The two available QPUs are Starmon-5 and Spin-2 which have five and two qubits respectively. These QPUs suffer from noise, limited coherence time, limited qubit connectivity, and each QPU has a specific allowed gate set. Note that these gate sets are universal, and non-native gates are decomposed into supported gates depending on the QPU. Regardless of the limitations of these QPUs, they are an essential part in research towards better quantum hardware and quantum algorithms.

In the absence of large-scale and fault-tolerant quantum computers, quantum computer simulation is critical for developing and testing quantum algorithms. Quantum computer simulators often use one of three classical algorithms to simulate quantum circuits. These algorithms have different trade-offs in the resources that they use. First, there is the naive Schrödinger simulation algorithm in which the complex amplitudes of a quantum state are represented in a state-vector, gates are represented as matrices, and state evolution is done by matrix-vector multiplication. This algorithm uses $O(2^n)$ space and time, where n is the amount of qubits [69]. Second, there is the Feynman

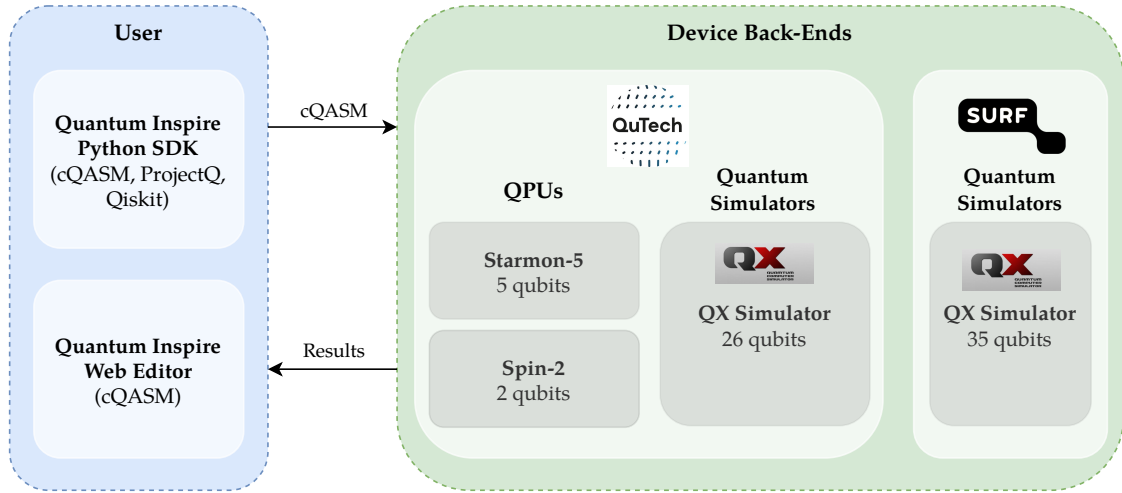


Figure 4.1: Overview of the Quantum Inspire workflow. Users can submit quantum circuits written in cQASM using the web editor or Python SDK. After the program has been run, the results are returned to the user. The quantum circuit can be executed on one of the QPUs or simulated using one of the QX simulator back-ends. QuTech’s simulator back-end supports simulations up to 26 qubits, while SURF’s simulator back-end supports simulations up to 35 qubits.

simulation algorithm which makes use of the path integral formulation of quantum mechanics [70]. This approach uses $O(m + n)$ space, but takes $O(4^m)$ time, where m is the amount of gates [69]. So while using a linear amount of space, the time it takes grows exponential with the number of gates. The number of gates even for small computations can reach the order of thousands. Supercomputers may be able to simulate quantum circuits with a thousand gates using the Schrödinger algorithm, but not even the most powerful supercomputer could simulate such circuit using the Feynman algorithm. On the other hand, the Feynman algorithm can simulate a quantum circuit with a few gates and 100 qubits, while the Schrödinger algorithm is greatly limited in the amount of qubits it can simulate due to its exponential space requirement. Finally, Aaronson and Chen [69] described an algorithm that offers best of both worlds: the Schrödinger-Feynman algorithm. This algorithm uses $O(d \cdot n \log n)$ space and $O(d^n)$ time, where d is the depth of the circuit. These different simulation algorithms allow you to trade time for space efficiency and vice versa, but all of them require exponential time and/or space, greatly limiting the size of simulations we can do.

Quantum Inspire offers two simulator back-ends which use QX [71] as quantum computer simulator. The QX simulator uses the Schrödinger simulation algorithm focused on efficiently using sparse operators. One of the simulator back-ends is hosted by QuTech on a commodity cloud-based server with 4 GB of memory, which can run simulations up to 26 qubits. The other simulator back-end is hosted on SURF’s computer cluster Lisa which consists of several hundreds of multi-core nodes. The largest node available on Lisa has 1.5 TB of memory, which supports simulations up to 35 qubits.

When quantum circuits are submitted to Quantum Inspire, they are handed to a job scheduler which will attempt to schedule the job as quickly as possible when the requested resources are available. For the QPUs and QuTech simulator back-ends, the jobs are simply placed in a queue which are executed in first-in-first-out order. The wait time for these back-ends ranges from seconds to minutes. For the SURF simulator back-end, the jobs are handled by a batch system. Because SURF’s infrastructure is

shared with many other users, wait times can vary from minutes to hours. Especially large simulations which require a large amount of resources may take a long time to schedule.

4.2 SURF High Performance Computing

SURF offers an integrated ICT research infrastructure and provides services in the areas of computing, data storage, visualization, networking, and cloud. This report is mainly interested in the computing services, and more precisely the Lisa and Cartesius cluster computers. These cluster computers can be thought of as a network of computers (nodes). Each node has its own central processing units (CPUs), memory, disk space, and a shared file system. Cluster computers are often used for high performance computing (HPC), where a large amount of resources (CPUs, graphics processing units (GPUs), memory, etc.) are used to solve computationally expensive problems. Lisa and Cartesius use the Slurm workload manager [72] as job scheduler. The specifications of the Lisa and Cartesius cluster computers are described in Table 4.1.

Resource	Lisa	Cartesius
CPUs	5304	47,776
Memory	41.7 TB	130 TB
GPUs	216	132
Flop/s (64-bit floating point)	367.9 Tflop/s	1.843 Pflop/s
Disk space	400 TB home	180 TB home, 7.7 PB project + scratch

Table 4.1: Specification of the Lisa and Cartesius cluster computers. Here Flop/s refers to floating point operations per second, or the theoretical peak performance. Cartesius is the Dutch national supercomputer meant for large computations that require many cores, large symmetric multiprocessing nodes, large amounts of memory, a fast interconnect, large amount of disk space, or a fast I/O subsystem. Lisa is meant for tasks that need large computing capacities, but do not need the facilities of a real supercomputer.

The Lisa and Cartesius cluster computers consist of login nodes and batch nodes. One connects to the system through the login node, which is an interactive system similar to a personal computer. These login nodes can be used for light tasks such as writing and compiling software, writing job scripts, and preparing data. From these login nodes, you interact with the batch nodes. Interacting with batch nodes works different than a regular computer: instead of submitting a command and immediately having it run, with cluster computers you submit a job script that runs as soon as sufficient resources are available. A job script is essentially a recipe of commands that you want your job to execute. Through this script you can load the required modules, handle input and output data processing, define the amount of resources your job requires, and handle the running of the actual application. The job script is submitted to the job scheduler and will be run on a batch node as soon as the requested resources are available.

4.3 Hybrid Quantum-Classical Workflows

When researching and running HQCAs, different problems require different workflows. For example, small HQCA experiments do not require HPC facilities, unless one wants to run the experiment many times in parallel with varying parameters. In Table 4.2

the five different workflows which are considered in this report are described. The order in which the workflows are displayed can be thought of as progressing from small experimentation and prototyping to (large) well-defined experiments. First, the Local-Local hybrid quantum-classical workflow is often enough for a lot of use cases as most research towards HQCAs is theoretical and does not require HPC facilities or access to a QPU. The Local-HPC workflow is useful for large quantum circuit simulations that exceed the resources of modern laptops. In such case the HPC-HPC workflow may also be used to reduce network overhead and repeated scheduler wait times. Generally, one wants the quantum and classical part to be as close to each other as possible, and ideally on the same network. This is usually possible for quantum simulators, but becomes more difficult when dealing with QPUs. The Local-QPU workflow is useful for experimenting and testing out quantum hardware. In the context of HQCAs, this can be useful for running small experiments and seeing how noise affects the performance of the algorithm. For larger and more well-defined experiments, the HPC-QPU workflow offers the advantage of being able to run experiments unattended. One could for example define multiple experiments with different parameters and run it over night on SURF's cluster computer. This is also the biggest advantage of the HPC-HPC workflow — one could run many experiments in parallel instead of being limited to local resources.

Workflow name	Classical part	Quantum part	Use case(s)
Local-Local	Local	Simulation (Local)	Small experimentation/prototyping
Local-HPC	Local	Simulation (HPC)	Large quantum simulations that can not be run locally
Local-QPU	Local	QPU	Small experimentation/exploration of hardware
HPC-QPU	HPC	QPU	Well-defined/multiple experiments, hybrid computation with expensive classical part
HPC-HPC	HPC	Simulation (HPC)	Multiple experiments, expensive classical and/or quantum part, high performance simulations

Table 4.2: Relevant hybrid quantum-classical workflows and their use cases. In the classical part column the platform on which the classical computation runs: locally on a laptop or personal computer, or on SURF's HPC center. Equivalently, the quantum part columns refers to the platform the quantum computation runs on: local simulation, QPU, or simulation on SURF's HPC center.

4.4 Hybrid Quantum-Classical Workflow Analysis

In this section we define and analyze the efficiency of executing HQCAs using the different workflows defined in the previous section. There are many variables to be considered when talking about the efficiency of a hybrid quantum-classical workflow:

network latency, circuit optimization, scheduler wait times, optimizer choice, etc. In the end, the aim is to reduce the amount of time it takes to run a full experiment. The variables we consider to describe the efficiency of a hybrid quantum-classical experiment in this report are as follows:

- $T_Q(n, d, s)$ describes the quantum circuit execution time where n is the number of qubits, d the depth of the quantum circuit, and s the number of shots.
- T_L describes the latency between the classical and quantum computation. This includes latencies such as network overhead, scheduler wait time, resource initialization, etc.
- $T_C(p)$ describes the time taken by the classical optimization algorithm for p parameters.

We then characterize the time taken by a hybrid quantum-classical experiment by the sum of these variables:

$$T_{\text{total}} = T_Q(n, d, s) + T_L + T_C(p), \quad (4.1)$$

and a lower value is considered more efficient.

To get an initial understanding on the efficiency of different quantum back-ends we implement and benchmark the QAOA to solve the maximum cut (Max-Cut) problem on a 2-regular graph with 4 vertices (qubits) and edges. The details of the Max-Cut problem and the QAOA implementation are described in Appendix A. This report will mainly focus on the quantities T_Q and T_L as they are the largest and most relevant. Some work related to optimizing the classical optimization part has been done by Lavrijsen, Tudor, Müller, *et al.* [73] and Sung, Harrigan, Rubin, *et al.* [74]. The benchmark results are displayed in Table 4.3. The efficiency of these different quantum parts should not be compared to each other as each quantum part is a different tool with a different application. Rather, it can be interpreted as a demonstration of how different quantum parts behave on the same problem.

Quantum part	$T_Q(n, d, s)$	T_L	$T_Q(n, d, s) + T_L$
Simulation (Local, Qiskit)	0.009s \pm 0.001s	0.010s \pm 0.022s	0.019s \pm 0.023s
Simulation (HPC)	0.016s \pm 0.001s	11.734s \pm 2.029s	11.751s \pm 2.029s
QPU (Starmon-5)	0.622s \pm 0.000s	27.742s \pm 0.610s	28.365s \pm 0.610s

Table 4.3: Benchmarks of the time taken by different quantum parts in the context of a hybrid quantum-classical computation. We measure the quantum execution time $T_Q(n, d, s)$ and latency T_L for 150 quantum circuit executions for every quantum part, with $n = 4$, $d = 32$, and $s = 4096$. The values shown are the mean and standard deviation of the 150 runs.

Given the measurements from Table 4.3, we can estimate how much time it would take to execute a meaningful and practical HQCA experiment. In the context of running the VQE on a hardware back-end, Kandala, Mezzacapo, Temme, *et al.* [48] estimate that with the current architecture in the order of 10^5 (best case) to 10^6 (worst case) quantum circuit executions are needed to achieve results within chemical accuracy. This is a very rough estimate, as the amount of circuit executions required depends on the problem at hand, desired accuracy, hardware error rates, choice of ansatz, etc. Using the total average time taken per circuit execution by the Starmon-5 QPU from Table 4.3, this would

take in the order of 32 to 328 days just in quantum execution and wait time. On top of that, there is an additional development time that should be considered, which includes debugging and testing quantum circuit behavior using quantum circuit simulators. In the following sections we will explore different methods of potentially increasing the efficiency of the execution of HQCAs.

4.4.1 Development Cycle

The development cycle of HQCAs usually consists of three phases of quantum circuit execution: noiseless simulation, noisy simulation, and execution on hardware. Quantum Inspire’s QPU and simulation back-ends are all in the cloud, introducing a network overhead. While this is hard to avoid for the QPU back-ends, for some applications of the simulator back-ends this network overhead may be avoided through local simulation. This is especially true for simulations of approximately 24 qubits and lower, for which the benefits of HPC do not yet become substantial. For this range of simulations, the network and possible job queue overhead greatly exceeds the actual quantum circuit execution time. To avoid this unnecessary overhead for small simulations, support for a local simulation back-end in the Quantum Inspire SDK can be useful. As per the measurements from Table 4.3, local simulation is a factor of ~ 1000 faster for 4-qubit simulations over HPC simulation because of the network and job scheduler overhead.

Another way to increase the efficiency of the HQCA development cycle is to reduce the amount of circuits that need to be executed on a QPU. This can be achieved if one has access to a noisy quantum simulator which models the target QPU’s hardware noise well. QX simulator supports the quantum depolarizing channel [30, Section 8.3.4], which is a standard model for noise in quantum systems. However, the quantum depolarizing channel does not accurately represent the errors encountered on actual QPUs. By adding support for better user-configurable error channels that correspond more to actual hardware noise, the amount of quantum circuits that need to be executed on actual QPUs can be decreased. One can then use the noisy quantum simulator for a larger part of the development process, which will increase the development efficiency as quantum circuit simulation is often more accessible and has less overhead cost.

4.4.2 QX Simulator Batch Parameters

When submitting a job to the job scheduler on SURF’s HPC center, one can define the resources required for the job through batch parameters. Choosing the right parameters for a job can optimize its execution time and help the job scheduler use the system’s resources as optimal as possible. In the context of the SURF QX simulator back-end, we are interested in figuring out the optimal number of cores and memory to use for quantum circuit simulations of different qubit counts. Choosing the optimal number of cores improves both the quantum simulation execution time (T_Q) and scheduling efficiency (T_L), whereas choosing the right amount of memory ensures that the simulation has access to enough memory and increases scheduling efficiency. We run all experiments and benchmarks on SURF’s cluster computer Lisa on the quantum computing nodes (`fat_soil_shared` partition). This partition consists of four nodes with 1.5 TB memory, two NVIDIA TitanRTX GPUs, and four sockets with 10 CPUs per socket. These quantum HPC nodes are especially tailored for the development and execution of (hybrid) quantum simulations.

Number of Cores Parameter

To determine the optimal number of cores to use for a given number of qubits, we measured the execution time of QX simulator for quantum circuits with different number of qubits. The quantum circuits used are circuits generated by our QAOA implementation on the Max-Cut problem. For each number of qubits n , we generate a QAOA Max-Cut circuit on a random 2-regular graph with n vertices and $p = 1$. These random graphs are generated with no self-loops or parallel edges according to the algorithm described by Steger and Wormald [75]. The respective quantum circuit for a 2-regular graph with n vertices has a depth of $6n + 2$. This means that the depth increases as the number of qubits and problem size increases, as is the case with most realistic implementations of quantum algorithms. The thread affinity was set up in a way that the threads run as much on the same CPU socket(s) as possible according to the layout from Table 4.4.

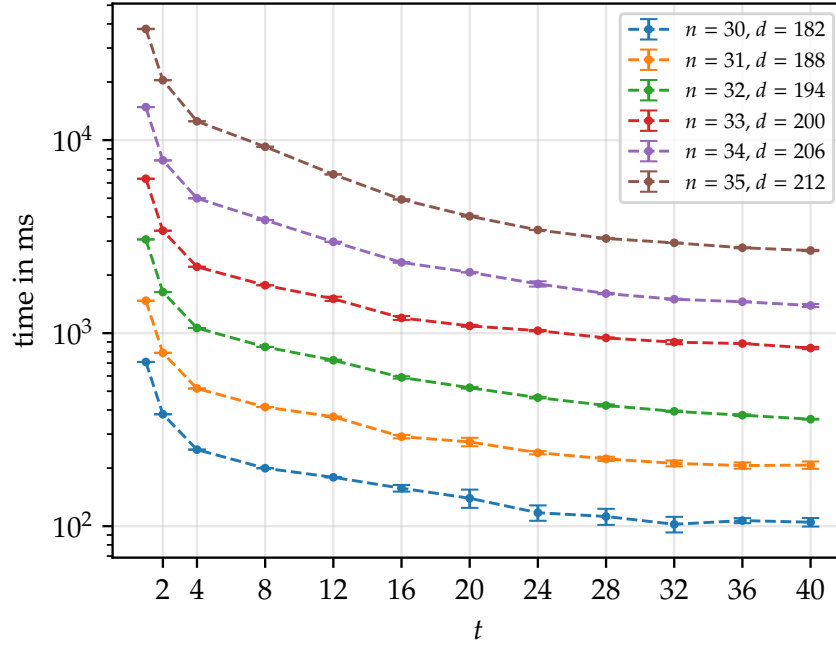
CPU socket	CPUs
0	0, 4, 8, 12, 16, 20, 24, 28, 32, 36
1	1, 5, 9, 13, 17, 21, 25, 29, 33, 37
2	2, 6, 10, 14, 18, 22, 26, 30, 34, 38
3	3, 7, 11, 15, 19, 23, 27, 31, 35, 39

Table 4.4: Overview of the layout of the sockets and CPUs on the `fat_soil_shared` quantum HPC nodes. Each node has 4 sockets with 10 CPUs per socket for a total of 40 CPUs.

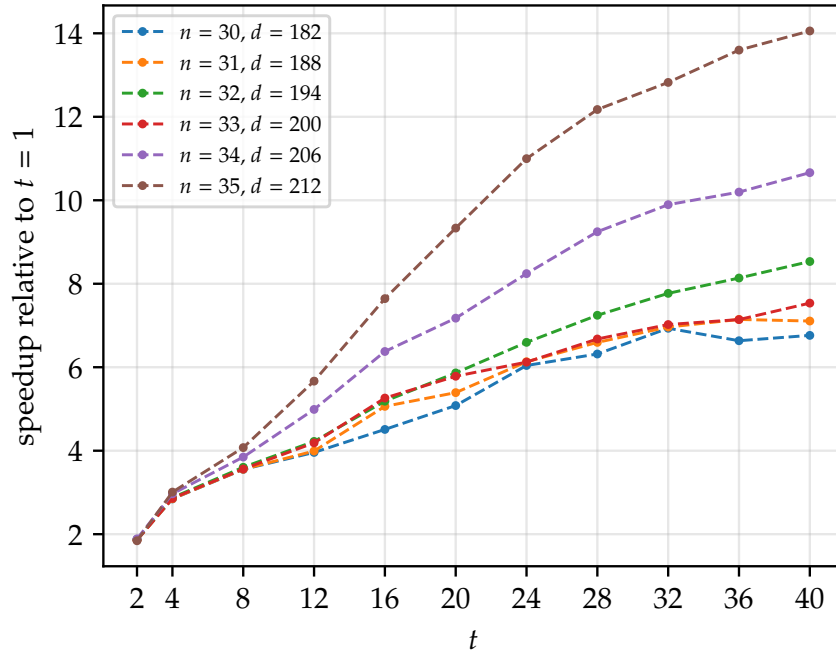
The benchmark results are displayed in Table 4.5. The timings and speedups of simulations with qubit counts ≥ 30 are visualized in Figure 4.2a and Figure 4.2b respectively. For lower qubit counts adding more cores does not or only slightly improve the execution time. Given the limited amount of resources available on the quantum HPC nodes and the minor speedup gained, simulations below 20 qubits should probably not be given more than one or two cores. For simulations of 20 qubits and above, increasing the amount of cores becomes more fruitful. As the number of qubits increases, the execution time increases exponentially, and for simulations of 20 qubits and more, increasing the amount of cores begins to save time in the order of seconds and even minutes. These higher qubit simulations can be sped up an order of 7 to 10 by using 32 or more cores. The speedup stagnates after 32 cores, and while for some cases the speedup continues to increase beyond 32 cores, it is a relatively small speedup. Depending on the seasonality (how many people are actively using the system), 32 cores per simulation might occupy too many resources. In this case, a speedup on the order of 4 to 6 may still be achieved with 16 cores. The trade-off to be considered when choosing the amount of cores to use per number of qubits depends on the available resources, seasonality, and speedup gained. The performance of quantum circuit simulation can be further pushed by using GPUs as shown in [76, Appendix B], however QX simulator does not currently support GPU simulation.

n	d	$t = 1$	$t = 2$	$t = 4$	$t = 8$	$t = 12$	$t = 16$	$t = 20$	$t = 24$	$t = 28$	$t = 32$	$t = 36$	$t = 40$
4	26	0.026s	0.023s	0.027s	0.027s	0.063s	0.027s	0.025s	0.027s	0.031s	0.031s	0.032s	0.043s
6	38	0.023s	0.024s	0.023s	0.027s	0.024s	0.028s	0.028s	0.028s	0.032s	0.032s	0.028s	0.031s
8	50	0.023s	0.025s	0.025s	0.028s	0.027s	0.028s	0.025s	0.027s	0.031s	0.030s	0.028s	0.035s
10	62	0.023s	0.024s	0.026s	0.025s	0.027s	0.030s	0.030s	0.028s	0.028s	0.028s	0.027s	0.041s
12	74	0.025s	0.026s	0.023s	0.028s	0.027s	0.029s	0.035s	0.029s	0.027s	0.034s	0.029s	0.045s
14	86	0.032s	0.032s	0.027s	0.030s	0.028s	0.030s	0.032s	0.028s	0.033s	0.027s	0.029s	0.044s
16	98	0.047s	0.041s	0.041s	0.039s	0.036s	0.028s	0.042s	0.039s	0.038s	0.030s	0.038s	0.042s
18	110	0.106s	0.070s	0.053s	0.051s	0.048s	0.050s	0.039s	0.049s	0.037s	0.039s	0.042s	0.043s
20	122	0.417s	0.242s	0.145s	0.105s	0.096s	0.078s	0.072s	0.083s	0.085s	0.069s	0.075s	0.075s
22	134	2.023s	1.102s	0.719s	0.553s	0.445s	0.353s	0.292s	0.264s	0.234s	0.219s	0.213s	0.213s
24	146	8.853s	4.800s	3.164s	2.533s	2.101s	1.685s	1.446s	1.297s	1.187s	1.087s	1.026s	0.983s
26	158	38.943s	20.991s	13.787s	11.045s	9.239s	7.390s	6.668s	5.734s	5.212s	4.805s	4.512s	4.318s
28	170	165.9s	89.3s	58.5s	46.9s	42.0s	36.2s	36.3s	33.3s	28.9s	25.0s	22.8s	27.5s
30	182	709.1s	380.0s	249.1s	199.6s	179.1s	157.2s	139.5s	117.4s	112.2s	102.3s	106.9s	104.8s
31	188	1472.5s	790.6s	516.8s	414.0s	368.9s	290.9s	273.1s	240.2s	223.3s	211.5s	206.1s	207.2s
32	194	3057.3s	1632.7s	1063.3s	848.5s	724.1s	589.5s	521.4s	463.5s	422.0s	393.5s	375.7s	358.3s
33	200	6307.4s	3394.2s	2206.7s	1770.6s	1506.3s	1198.8s	1090.5s	1030.3s	944.7s	898.1s	883.2s	837.0s
34	206	14827.4s	7855.2s	4996.7s	3853.7s	2971.9s	2325.1s	2066.3s	1798.9s	1603.6s	1498.6s	1454.3s	1390.7s
35	212	37658.7s	20438.4s	12516.8s	9242.1s	6645.5s	4926.8s	4034.7s	3424.5s	3093.3s	2936.7s	2769.5s	2679.0s

Table 4.5: Benchmark results of the execution time of quantum circuits using QX simulator for different qubit and core counts. Here n is the number of qubits, d is the depth of the circuit, and t is the number of cores. Execution times shown are the average of 10 repeated circuit executions. The optimal number of cores t is highlighted for every n .



(a) Time taken by QX simulator on the QAOA circuit for different number of threads t . The standard deviation is plotted as error bars.



(b) Speedup factor gained by using t cores relative to single core ($t = 1$) simulation.

Figure 4.2: Visualizations of the QX simulator benchmark results for $n \geq 30$. Figure 4.2a shows the average total time taken by QX simulator for an n -qubit simulation with t cores, and Figure 4.2b plots the speedup factor gained by adding more cores for different qubit counts.

Memory Parameter

The `fat_soil_shared` partition is composed of shared nodes, which means that multiple users can run jobs on the same node and resources like memory and disk space are shared. For this reason, to ensure that a simulation has access to enough memory, the simulation job should specify the amount of memory it needs. By default, a job on the `fat_soil_shared` partition is given 37.5 GB of memory per requested core. This is not a fit default for QX simulator, as its memory usage grows with the number of qubits, not necessarily the amount of cores it uses. To ensure that simulations request a reasonable amount of memory, we measure the memory usage of QX simulator running random circuits of depth 2 for different number of qubits. This is done by spawning a QX simulator process and polling the operating system every 0.1 millisecond for the respective process' resident set size (RSS) memory usage. We then take the highest amount of memory measured by polling and take that as the "total" memory usage of the simulation. This is repeated 10 times per n qubits. The results are displayed in Table 4.6. We plot the data in Figure 4.3, along with a fitted exponential function to predict the amount of memory \hat{m} required (in GB) for a simulation of n qubits:

$$\hat{m}(n) = (3.20017417 \cdot 10^{-8}) \cdot e^{0.693145626n} + 0.0210389214. \quad (4.2)$$

This function is especially useful to define the memory requirement for simulations of $n \geq 31$ qubits, where the default of 37.5 GB of memory per requested core becomes too small.

n	RSS memory usage	n	RSS memory usage
1	18.96 MB \pm 918.1 kB	19	36.45 MB \pm 715.7 kB
2	19.72 MB \pm 1.8 MB	20	53.87 MB \pm 915.5 kB
3	19.97 MB \pm 971.3 kB	21	88.12 MB \pm 159.4 kB
4	19.72 MB \pm 1.1 MB	22	153.70 MB \pm 1.6 MB
5	20.28 MB \pm 1.7 MB	23	288.00 MB \pm 858.9 kB
6	19.52 MB \pm 1.0 MB	24	555.85 MB \pm 742.7 kB
7	20.74 MB \pm 2.2 MB	25	1.09 GB \pm 2.5 MB
8	19.89 MB \pm 1.5 MB	26	2.17 GB \pm 892.8 kB
9	19.79 MB \pm 1.1 MB	27	4.32 GB \pm 963.7 kB
10	19.87 MB \pm 991.8 kB	28	8.61 GB \pm 1.0 MB
11	19.45 MB \pm 1.0 MB	29	17.20 GB \pm 1.3 MB
12	19.48 MB \pm 1.2 MB	30	34.38 GB \pm 2.3 MB
13	20.00 MB \pm 1.9 MB	31	68.74 GB \pm 781.9 kB
14	19.64 MB \pm 1.0 MB	32	137.46 GB \pm 1.0 MB
15	20.77 MB \pm 725.0 kB	33	274.90 GB \pm 901.2 kB
16	23.35 MB \pm 2.2 MB	34	549.78 GB \pm 1.2 MB
17	21.51 MB \pm 935.2 kB	35	1.10 TB \pm 995.4 kB
18	26.94 MB \pm 1.5 MB		

Table 4.6: QX simulator RSS memory usage for varying number of qubits. The measurements shown are the average taken of 10 runs.

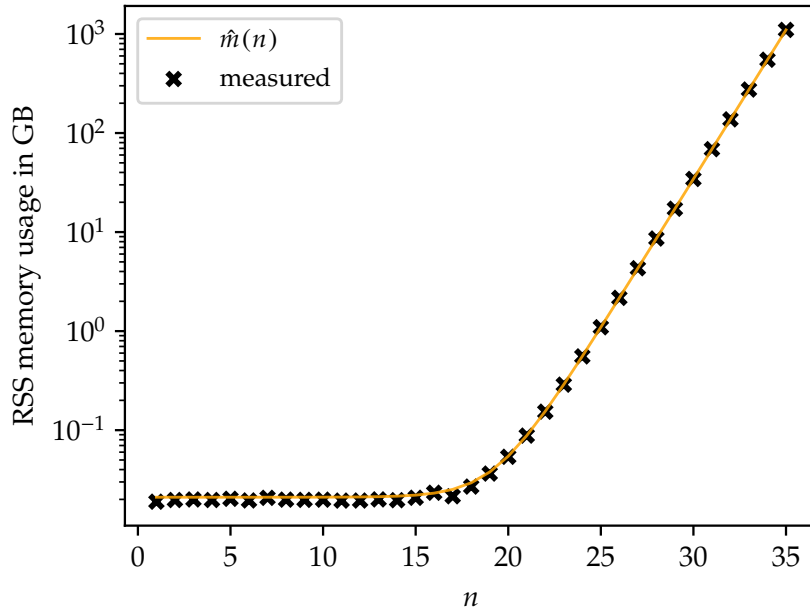


Figure 4.3: Plot of the RSS memory usage of QX simulator as the number of qubits increases. We use non-linear least squares [77] to fit an exponential function $\hat{m}(n) = a \cdot e^{bn} + c$ to the data, which gives $a = 3.20017417 \cdot 10^{-8}$, $b = 0.693145626$, and $c = 0.0210389214$ with $R^2 = 0.99$. We fit the function to the mean of 10 runs plus its standard deviation to err on the high side to ensure that the simulation has more than enough memory available.

4.4.3 Parallel Quantum Circuit Execution

In the VQE one uses a quantum computer to measure the energy of a molecular Hamiltonian. As described in Section 3.2, this is done by decomposing the Hamiltonian into a linear combination of Pauli operators (Equation 3.10). Consider the following simple one-qubit Hamiltonian:

$$H = g_0 Z_0 + g_1 X_0, \quad (4.3)$$

for which the energy is given as

$$\langle H \rangle = g_0 \langle Z_0 \rangle + g_1 \langle X_0 \rangle, \quad (4.4)$$

where g_0 and g_1 are some constants. Since the Pauli terms $\langle Z_0 \rangle$ and $\langle X_0 \rangle$ do not commute and measurement is destructive, we need to run the same quantum circuit twice to measure in different directions: once for measuring $\langle Z_0 \rangle$, and once for measuring $\langle X_0 \rangle$. These circuit evaluations are independent of each other and could be executed in parallel to cut the calculation time in half. The number of measurement directions required per energy measurement increases rapidly for molecular Hamiltonians: the LiH and BeH₂ molecules require 25 and 44 measurement directions respectively [48, Supplementary Information, Section III]. For simulator back-ends this parallelization is easily achieved by simply running multiple simulations in parallel. Even more, because we have full access to the state vector in quantum simulation, one can simply save the final state vector of the quantum circuit before measurement and measure in any number of ways you like. For QPU back-ends, this kind of parallelization requires access to multiple QPUs. At this moment, most development in quantum hardware is focused on the quality, and not the quantity of QPUs. However, access to multiple NISQ devices may be required

to make the execution of HQCAs feasible and practical. Assuming even access to two QPUs, the estimate of the execution time of a meaningful HQCA experiment given in Section 4.4 can be approximately halved from ranging from 32 to 328 days, to ranging from 16 to 164 days.

4.4.4 Quantum Circuit Hardware Batching

For the QPU back-ends of Quantum Inspire, most time is spent on circuit compilation and uploading to control hardware (T_L in Table 4.3). The control hardware for Quantum Inspire’s QPU back-ends support uploading multiple quantum circuits in one batch, which can be used to reduce the overall amount of time spent on these overheads. The QPUs can still only execute one circuit at a time, but by uploading multiple circuits in the same batch, there is only a overhead for the first circuit. In this section, we will look at how we can take advantage of this circuit batching to reduce the time taken to execute a HQCA experiment on Quantum Inspire’s QPU back-ends.

Gradient Circuits Batching

Recall from Section 3.1.3 that we can analytically calculate the gradient of a quantum circuit through the parameter-shift rule. That is, given a parameterized quantum gate of the form $G(\theta) = e^{-i\theta/2H}$, one can calculate the derivative of G with respect to θ by executing the quantum circuit twice with shifted parameters $\theta \pm \pi/2$ (Equation 3.5). For a circuit with m parameterized gates, $2m$ quantum circuit executions are required to calculate the gradient. These gradient circuits can be generated in parallel and submitted in a single batch.

Particle Swarm Optimization

Gradient based optimization algorithms can make use of circuit batching as described in the previous section, but non-gradient methods are generally limited to executing circuits sequentially. The iterative optimization algorithms discussed in this report so far consider one candidate solution in the parameter space at the same time. The particle swarm optimization (PSO) method optimizes a function by considering a population of candidate solutions (particles), and moving these particles around according to the particle’s position and velocity [78, 79]. A particle moves through the search space based on its local best known position, but is also influenced by positions found by other particles.

4.4.5 Summary

This chapter described different methods for improving the efficiency of the execution of HQCAs using the Quantum Inspire quantum computing platform and SURF’s HPC center. To summarize, we propose the following changes:

- In Section 4.4.1 we show that for quantum circuit simulations that do not require HPC facilities, local simulation is a factor of 1000 faster than running simulations in the cloud. In addition, by adding support for error channels that correspond more to actual hardware noise in QX simulator, one can reduce the amount of quantum circuits that need to be executed on actual QPUs, reducing the overall time of a HQCA experiment.

- Section 4.4.2 analyzes the performance of QX simulator on SURF’s cluster computer Lisa. Based on benchmark data we propose methods for determining the optimal batch parameters for QX simulations. This allows for more efficient scheduling and improves the simulation execution speed up to a factor of 14.
- Section 4.4.3 shows how access to multiple QPUs can improve the execution time of VQE implementations by parallelizing measurements in different directions.
- Section 4.4.4 describes how one can make use of the control hardware’s ability to upload multiple quantum circuits in one batch to reduce the total time spent on overheads such as circuit compilation and uploading to control hardware.

5

Conclusion

Bibliography

- [1] P. Benioff, "The computer as a physical system: A microscopic quantum mechanical hamiltonian model of computers as represented by turing machines", *Journal of statistical physics*, vol. 22, no. 5, pp. 563–591, 1980.
- [2] I. I. Manin, *Vychislimoe i nevychislimoe*. Sov. radio, 1980.
- [3] R. P. Feynman, "Simulating physics with computers", *Int. J. Theor. Phys*, vol. 21, no. 6/7, 1982.
- [4] P. W. Shor, "Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer", *SIAM review*, vol. 41, no. 2, pp. 303–332, 1999.
- [5] L. K. Grover, "A fast quantum mechanical algorithm for database search", in *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, 1996, pp. 212–219.
- [6] S. McArdle, S. Endo, A. Aspuru-Guzik, *et al.* (2018). Quantum computational chemistry. arXiv: 1808.10402 [quant-ph].
- [7] C. H. Bennett and G. Brassard, "Quantum cryptography", *Theoretical Computer Science*, vol. 560, no. P1, pp. 7–11, 2014.
- [8] J. Biamonte, P. Wittek, N. Pancotti, *et al.*, "Quantum machine learning", *Nature*, vol. 549, no. 7671, pp. 195–202, 2017.
- [9] A. G. Fowler, M. Mariantoni, J. M. Martinis, *et al.*, "Surface codes: Towards practical large-scale quantum computation", *Physical Review A*, vol. 86, no. 3, p. 032 324, 2012.
- [10] J. Preskill, "Quantum computing in the NISQ era and beyond", *Quantum*, vol. 2, p. 79, 2018.
- [11] S. Bravyi, D. Gosset, R. Koenig, *et al.*, "Quantum advantage with noisy shallow circuits", *Nature Physics*, vol. 16, no. 10, pp. 1040–1045, 2020.
- [12] S. Endo, Z. Cai, S. C. Benjamin, *et al.*, "Hybrid quantum-classical algorithms and quantum error mitigation", *Journal of the Physical Society of Japan*, vol. 90, no. 3, p. 032 001, 2021.
- [13] QuTech. (2018). Quantum inspire home, [Online]. Available: <https://www.quantum-inspire.com/> (visited on 2021-02-24).
- [14] SURF. (2021). SURFsara systems, [Online]. Available: <https://userinfo.surfsara.nl/systems/> (visited on 2021-02-24).
- [15] A. Church *et al.*, "A note on the entscheidungsproblem", *J. Symb. Log.*, vol. 1, no. 1, pp. 40–41, 1936.
- [16] A. M. Turing, "On computable numbers, with an application to the entscheidungsproblem", *Proceedings of the London mathematical society*, vol. 2, no. 1, pp. 230–265, 1937.
- [17] T. H. Cormen, C. E. Leiserson, R. L. Rivest, *et al.*, *Introduction to algorithms*. MIT press, 2009.

- [18] S. Arora and B. Barak, *Computational complexity: A modern approach*. Cambridge University Press, 2009.
- [19] J. Hartmanis and R. E. Stearns, "On the computational complexity of algorithms", *Transactions of the American Mathematical Society*, vol. 117, pp. 285–306, 1965.
- [20] N. Dershowitz and Y. Gurevich, "A natural axiomatization of computability and proof of church's thesis", *Bulletin of symbolic logic*, vol. 14, no. 3, pp. 299–350, 2008.
- [21] C. H. Bennett, "Logical reversibility of computation", *IBM journal of Research and Development*, vol. 17, no. 6, pp. 525–532, 1973.
- [22] D. Deutsch, "Quantum theory, the church-turing principle and the universal quantum computer", *Proceedings of the Royal Society of London. A. Mathematical and Physical Sciences*, vol. 400, no. 1818, pp. 97–117, 1985.
- [23] F. Arute, K. Arya, R. Babbush, *et al.*, "Quantum supremacy using a programmable superconducting processor", *Nature*, vol. 574, no. 7779, pp. 505–510, 2019.
- [24] S. Aaronson, *Quantum computing since Democritus*. Cambridge University Press, 2013.
- [25] O. Goldreich, "In a world of $P = BPP$ ", in *Studies in Complexity and Cryptography. Miscellanea on the Interplay between Randomness and Computation*, Springer, 2011, pp. 191–232.
- [26] N. Nisan and A. Wigderson, "Hardness vs randomness", *Journal of computer and System Sciences*, vol. 49, no. 2, pp. 149–167, 1994.
- [27] E. Bernstein and U. Vazirani, "Quantum complexity theory", *SIAM Journal on computing*, vol. 26, no. 5, pp. 1411–1473, 1997.
- [28] S. Aaronson, "BQP and the polynomial hierarchy", in *Proceedings of the forty-second ACM symposium on Theory of computing*, 2010, pp. 141–150.
- [29] G. M. D'Ariano, M. G. Paris, and M. F. Sacchi, "Quantum tomography", *Advances in Imaging and Electron Physics*, vol. 128, pp. 206–309, 2003.
- [30] M. A. Nielsen and I. Chuang, *Quantum computation and quantum information*, 2002.
- [31] C. M. Dawson and M. A. Nielsen, "The Solovay-Kitaev algorithm", 2005. arXiv: quant-ph/0505030.
- [32] Y. Shi, "Both toffoli and controlled-not need little help to do universal quantum computation", 2002. arXiv: quant-ph/0205115.
- [33] A. W. Harrow, A. Hassidim, and S. Lloyd, "Quantum algorithm for linear systems of equations", *Physical review letters*, vol. 103, no. 15, p. 150 502, 2009.
- [34] C. Zalka, "Grover's quantum searching algorithm is optimal", *Physical Review A*, vol. 60, no. 4, p. 2746, 1999.
- [35] S. Sim, P. D. Johnson, and A. Aspuru-Guzik, "Expressibility and entangling capability of parameterized quantum circuits for hybrid quantum-classical algorithms", *Advanced Quantum Technologies*, vol. 2, no. 12, p. 1900 070, 2019.
- [36] M. Cerezo, A. Arrasmith, R. Babbush, *et al.*, "Variational quantum algorithms", 2020. arXiv: 2012.09265 [quant-ph].
- [37] M. Schuld, V. Bergholm, C. Gogolin, *et al.*, "Evaluating analytic gradients on quantum hardware", *Physical Review A*, vol. 99, no. 3, p. 032 331, 2019.

-
- [38] A. Mari, T. R. Bromley, and N. Killoran, "Estimating the gradient and higher-order derivatives on quantum hardware", *Physical Review A*, vol. 103, no. 1, p. 012 405, 2021.
- [39] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization", 2014. arXiv: 1412.6980 [cs].
- [40] L. Page, S. Brin, R. Motwani, *et al.*, "The pagerank citation ranking: Bringing order to the web.", Stanford InfoLab, Tech. Rep., 1999.
- [41] G. H. Golub and H. A. van der Vorst, "Eigenvalue computation in the 20th century", *J. Comput. Appl. Math.*, vol. 123, no. 1-2, pp. 35–65, 2000-11, ISSN: 0377-0427. DOI: 10.1016/S0377-0427(00)00413-1.
- [42] D. S. Abrams and S. Lloyd, "Quantum algorithm providing exponential speed increase for finding eigenvalues and eigenvectors", *Physical Review Letters*, vol. 83, no. 24, p. 5162, 1999.
- [43] A. Peruzzo, J. McClean, P. Shadbolt, *et al.*, "A variational eigenvalue solver on a photonic quantum processor", *Nature communications*, vol. 5, no. 1, pp. 1–7, 2014.
- [44] W. Ritz, "Über eine neue methode zur lösung gewisser variationsprobleme der mathematischen physik.", *Journal für die reine und angewandte Mathematik*, vol. 1909, no. 135, pp. 1–61, 1909.
- [45] A. D. Bookatz, "QMA-complete problems", 2012. arXiv: 1212.6312 [quant-ph].
- [46] J. Kempe, A. Kitaev, and O. Regev, "The complexity of the local hamiltonian problem", *SIAM Journal on Computing*, vol. 35, no. 5, pp. 1070–1097, 2006.
- [47] G. Ortiz, J. E. Gubernatis, E. Knill, *et al.*, "Quantum algorithms for fermionic simulations", *Physical Review A*, vol. 64, no. 2, p. 022 319, 2001.
- [48] A. Kandala, A. Mezzacapo, K. Temme, *et al.*, "Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets", *Nature*, vol. 549, no. 7671, pp. 242–246, 2017.
- [49] M. Born and R. Oppenheimer, "Zur quantentheorie der molekeln", *Annalen der physik*, vol. 389, no. 20, pp. 457–484, 1927.
- [50] P. J. O'Malley, R. Babbush, I. D. Kivlichan, *et al.*, "Scalable quantum simulation of molecular energies", *Physical Review X*, vol. 6, no. 3, p. 031 007, 2016.
- [51] I. Kassal, J. D. Whitfield, A. Perdomo-Ortiz, *et al.*, "Simulating chemistry using quantum computers", *Annual review of physical chemistry*, vol. 62, pp. 185–207, 2011.
- [52] R. Somma, G. Ortiz, J. E. Gubernatis, *et al.*, "Simulating physical phenomena by quantum networks", *Physical Review A*, vol. 65, no. 4, p. 042 323, 2002.
- [53] J. T. Seeley, M. J. Richard, and P. J. Love, "The bravyi-kitaev transformation for quantum computation of electronic structure", *The Journal of chemical physics*, vol. 137, no. 22, p. 224 109, 2012.
- [54] I. Shavitt and R. J. Bartlett, *Many-body methods in chemistry and physics: MBPT and coupled-cluster theory*. Cambridge university press, 2009.
- [55] A. G. Taube and R. J. Bartlett, "New perspectives on unitary coupled-cluster theory", *International journal of quantum chemistry*, vol. 106, no. 15, pp. 3393–3401, 2006.

-
- [56] B. Cooper and P. J. Knowles, "Benchmark studies of variational, unitary and extended coupled cluster methods", *The Journal of chemical physics*, vol. 133, no. 23, p. 234 102, 2010.
- [57] F. A. Evangelista, "Alternative single-reference coupled cluster approaches for multireference problems: The simpler, the better", *The Journal of chemical physics*, vol. 134, no. 22, p. 224 102, 2011.
- [58] D. Wecker, M. B. Hastings, and M. Troyer, "Progress towards practical quantum variational algorithms", *Physical Review A*, vol. 92, no. 4, p. 042 303, 2015.
- [59] J. Romero, R. Babbush, J. R. McClean, *et al.*, "Strategies for quantum computing molecular energies using the unitary coupled cluster ansatz", *Quantum Science and Technology*, vol. 4, no. 1, p. 014 008, 2018.
- [60] H. R. Grimsley, S. E. Economou, E. Barnes, *et al.*, "An adaptive variational algorithm for exact molecular simulations on a quantum computer", *Nature communications*, vol. 10, no. 1, pp. 1–9, 2019.
- [61] E. Farhi, J. Goldstone, and S. Gutmann, "A quantum approximate optimization algorithm", 2014. arXiv: 1411.4028 [quant-ph].
- [62] A. Schrijver, *Combinatorial optimization: Polyhedra and efficiency*. Springer Science & Business Media, 2003, vol. 24.
- [63] E. Farhi and A. W. Harrow, "Quantum supremacy through the quantum approximate optimization algorithm", 2016. arXiv: 1602.07674 [quant-ph].
- [64] E. Farhi, J. Goldstone, S. Gutmann, *et al.*, "Quantum algorithms for fixed qubit architectures", 2017. arXiv: 1703.06199 [quant-ph].
- [65] T. Last, N. Samkharadze, P. Eendebak, *et al.*, "Quantum inspire: QuTech's platform for co-development and collaboration in quantum computing", in *Novel Patterning Technologies for Semiconductors, MEMS/NEMS and MOEMS 2020*, International Society for Optics and Photonics, vol. 11324, 2020, 113240J.
- [66] N. Khammassi, G. G. Guerreschi, I. Ashraf, *et al.*, "cQASM v1.0: Towards a common quantum assembly language", 2018. arXiv: 1805.09607 [quant-ph].
- [67] D. S. Steiger, T. Häner, and M. Troyer, "Projectq: An open source software framework for quantum computing", *Quantum*, vol. 2, p. 49, 2018.
- [68] G. Aleksandrowicz, T. Alexander, P. Barkoutsos, *et al.*, *Qiskit: an open-source framework for quantum computing*, version 0.7.2, 2019-01. DOI: 10.5281/zenodo.2562111. [Online]. Available: <https://doi.org/10.5281/zenodo.2562111>.
- [69] S. Aaronson and L. Chen, "Complexity-theoretic foundations of quantum supremacy experiments", in *32nd Computational Complexity Conference (CCC 2017)*, ser. Leibniz International Proceedings in Informatics (LIPIcs), vol. 79, Dagstuhl, Germany: Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2017, 22:1–22:67, ISBN: 978-3-95977-040-8. DOI: 10.4230/LIPIcs.CCC.2017.22.
- [70] R. P. Feynman, "Space-time approach to non-relativistic quantum mechanics", *Feynman's Thesis—A New Approach To Quantum Theory*, pp. 71–109, 2005.
- [71] N. Khammassi, I. Ashraf, X. Fu, *et al.*, "Qx: A high-performance quantum computer simulation platform", in *Design, Automation & Test in Europe Conference & Exhibition (DATE), 2017*, IEEE, 2017, pp. 464–469.

- [72] A. B. Yoo, M. A. Jette, and M. Grondona, "Slurm: Simple linux utility for resource management", in *Workshop on job scheduling strategies for parallel processing*, Springer, 2003, pp. 44–60.
- [73] W. Lavrijsen, A. Tudor, J. Müller, *et al.*, "Classical optimizers for noisy intermediate-scale quantum devices", in *2020 IEEE International Conference on Quantum Computing and Engineering (QCE)*, IEEE, 2020, pp. 267–277.
- [74] K. J. Sung, M. P. Harrigan, N. C. Rubin, *et al.*, "An exploration of practical optimizers for variational quantum algorithms on superconducting qubit processors", 2020. arXiv: 2005.11011 [quant-ph].
- [75] A. Steger and N. C. Wormald, "Generating random regular graphs quickly", *Combinatorics, Probability and Computing*, vol. 8, no. 04, pp. 377–396, 1999.
- [76] S. Oud. (2019). Simulation of quantum algorithms for solving machine learning problems, [Online]. Available: <https://github.com/soudy/quantum-ml-research-paper/blob/master/quantum-ml-research-paper.pdf>.
- [77] M. A. Branch, T. F. Coleman, and Y. Li, "A subspace, interior, and conjugate gradient method for large-scale bound-constrained minimization problems", *SIAM Journal on Scientific Computing*, vol. 21, no. 1, pp. 1–23, 1999.
- [78] J. Kennedy and R. Eberhart, "Particle swarm optimization", in *Proceedings of ICNN'95-international conference on neural networks*, IEEE, vol. 4, 1995, pp. 1942–1948.
- [79] Y. Shi and R. Eberhart, "A modified particle swarm optimizer", in *1998 IEEE international conference on evolutionary computation proceedings. IEEE world congress on computational intelligence (Cat. No. 98TH8360)*, IEEE, 1998, pp. 69–73.
- [80] P. Virtanen, R. Gommers, T. E. Oliphant, *et al.*, "Scipy 1.0: Fundamental algorithms for scientific computing in python", *Nature methods*, vol. 17, no. 3, pp. 261–272, 2020.

Glossary

COBYLA Constrained Optimization by Linear Approximations

CPU Central Processing Unit

cQASM Common Quantum Assembly Language

DFT Discrete Fourier Transform

GPU Graphics Processing Unit

HPC High Performance Computing

HQCA Hybrid Quantum-Classical Algorithm

Max-Cut Maximum Cut

NISQ Noisy Intermediate-Scale Quantum

PSO Particle Swarm Optimization

QAOA Quantum Approximate Optimization Algorithm

QFT Quantum Fourier Transform

QPEA Quantum Phase Estimation Algorithm

QPU Quantum Processing Unit

RSS Resident Set Size

SDK Software Development Kit

UCC Unitary Coupled Cluster

VQE Variational Quantum Eigensolver

Appendices

A

QAOA for Max-Cut

The aim of the Max-Cut problem is to find a set of vertices that maximizes the sum of the weights of the edges that are cut. It is a widely studied problem known to be NP-hard. Formally, given an undirected graph $G = (V, E)$ with $n = |V|$ vertices and non-negative weights $w_{j,k} = w_{k,j}$ on the edges $(j, k) \in E$, one looks to bipartition V into two sets S and $\bar{S} = V \setminus S$ such that the objective function L is maximized:

$$L(z) = \sum_{(j,k) \in E} w_{j,k} [z_j (1 - z_k) + z_k (1 - z_j)]. \quad (\text{A.1})$$

Here $z \in \{0, 1\}^n$ is a bit string that describes the bipartition as follows: if a vertex j is in partition S , then $z_j = 0$, and if a vertex j is in partition \bar{S} then $z_j = 1$. For example, Figure A.1 shows an example of a maximum cut on a graph with five vertices. The partition of the visualized cut can be described as the bit string $z = 00101$.

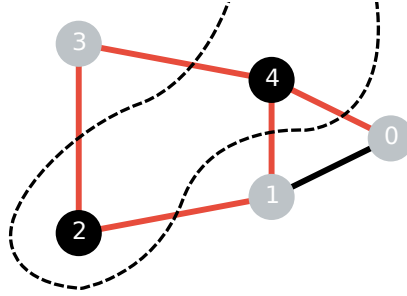


Figure A.1: Max-Cut example on a graph with five vertices and unit weights. The vertices are partitioned into two sets visualized as black and gray. The cut shown is a maximum cut with $L(z) = 5$, which can be thought of as the number of edges cut (shown in red).

The QAOA can be used for solving Max-Cut problems by assigning a vertex $j \in V$ to a qubit $|q_j\rangle$. A qubit $|q_j\rangle$ is in state $|0\rangle$ if a vertex j is in partition S , and state $|1\rangle$ if vertex j is in partition \bar{S} . To encode the objective function from Equation A.1, note that the objective function can be rewritten as follows:

$$L(z) = \frac{1}{2} \sum_{(j,k) \in E} w_{j,k} (1 - z_j z_k), \quad (\text{A.2})$$

where $z_j \in \{-1, 1\}$ for $j \in V$. This objective function can be represented by the following problem Hamiltonian:

$$H_L = \frac{1}{2} \sum_{(j,k) \in E} w_{j,k} (I - Z^{(j)} Z^{(k)}). \quad (\text{A.3})$$

This gives the problem unitary

$$U_L(\gamma) = e^{-i\gamma H_L} = \prod_{(j,k) \in E} e^{-i\gamma w_{j,k} (I - Z^{(j)} Z^{(k)})/2}, \quad (\text{A.4})$$

and the standard mixer unitary as defined in Section 3.3:

$$U_B(\beta) = e^{-i\beta H_B} = \prod_{j \in V} e^{-i\beta X^{(j)}}. \quad (\text{A.5})$$

A.1 Implementation

The QAOA was implemented in Python using the Project Q [67], Quantum Inspire [13], and SciPy [80] frameworks to solve the Max-Cut problem on the graph from Figure A.2. The optimal solutions for this graph are $z = 0101$ and $z = 1010$ with $L(z) = 4$. A

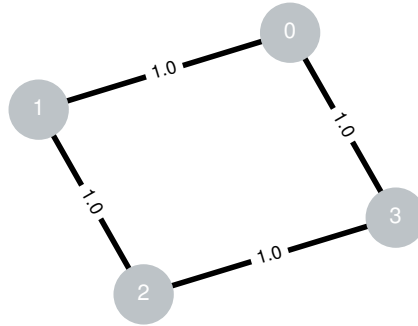
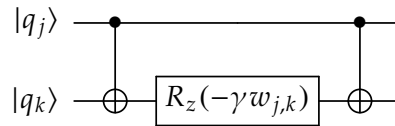


Figure A.2: Undirected 2-regular graph $G = (V, E)$ with $n = 4$ vertices $V = \{0, 1, 2, 3\}$ and 4 edges $E = \{(0, 1), (0, 3), (1, 2), (2, 3)\}$ with unit weight $w_{j,k} = w_{k,j} = 1$.

ZZ interaction $e^{-i\gamma w_{j,k} (I - Z^{(j)} Z^{(k)})/2}$ from the problem unitary $U_L(\gamma)$ (Equation A.4) is implemented as follows:



The X interaction $e^{-i\beta X^{(j)}}$ from the mixer unitary $U_B(\beta)$ (Equation A.5) is implemented as a R_x gate:

$$|q_j\rangle \text{ --- } [R_x(2\beta)] \text{ ---}$$

The complete quantum circuit for the QAOA for Max-Cut on this graph is shown in Figure A.3 and the respective cQASM for $p = 1$ is shown in Figure A.4.

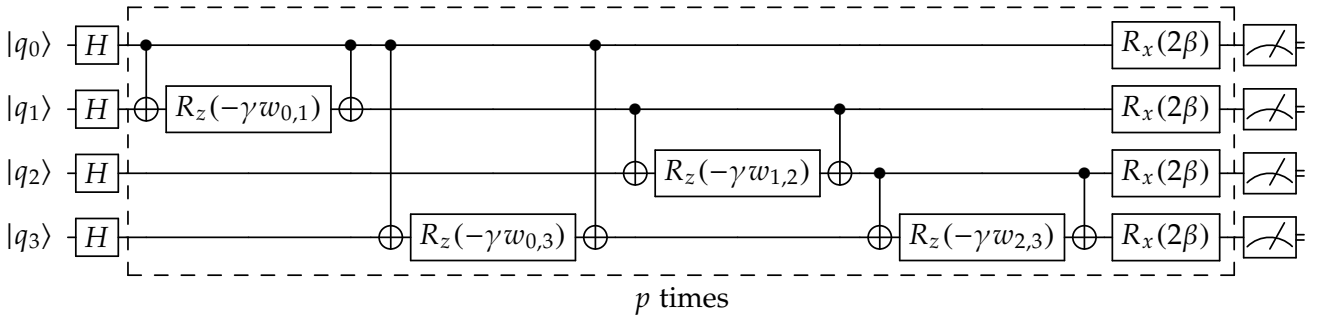


Figure A.3: Quantum circuit for the p -layer QAOA on the graph from Figure A.2. Each vertex $j \in V$ is represented by qubit $|q_j\rangle$. The circuit starts by preparing an equal superposition state, after which the problem unitary $U_L(\gamma)$ and mixer unitary $U_B(\beta)$ are applied p times. In general, the depth of the circuit is $1 + p(3m + 1)$, where p is the number of layers and m is the number of edges.

```

1  version 1.0
2
3  qubits 5
4
5  H q[0]
6  H q[1]
7  H q[2]
8  H q[3]
9  SWAP q[2], q[1]
10 CNOT q[3], q[2]
11 Rz q[2], 3.902484
12 CNOT q[3], q[2]
13 CNOT q[2], q[1]
14 Rz q[1], 3.902484
15 SWAP q[2], q[3]
16 CNOT q[2], q[0]
17 Rz q[0], 3.902484
18 SWAP q[2], q[1]
19 CNOT q[3], q[2]
20 SWAP q[0], q[2]
21 CNOT q[1], q[2]
22 CNOT q[0], q[2]
23 Rz q[2], 3.902484
24 CNOT q[0], q[2]
25 Rx q[0], 5.156352
26 Rx q[1], 5.156352
27 Rx q[2], 5.156352
28 Rx q[3], 5.156352
    
```

Figure A.4: cQASM for the 1-layer QAOA on the graph from Figure A.2 with $\gamma = 3.902484$ and $\beta = 5.156352$. Swap gates are added to deal with the limited qubit connectivity of the Starmon-5 hardware back-end.

The quantum circuit is implemented using the Qiskit quantum computing library. To find the optimal parameters $\gamma_{\text{opt}}, \beta_{\text{opt}}$ a classical optimizer provided by SciPy is used. The relevant cost function is

$$C_p(\gamma, \beta) = \langle \gamma, \beta | H_L | \gamma, \beta \rangle, \quad (\text{A.6})$$

and we look to solve the optimization problem

$$\gamma_{\text{opt}}, \beta_{\text{opt}} = \arg \max_{\gamma, \beta} C_p(\gamma, \beta). \quad (\text{A.7})$$

We then prepare the state $|\gamma_{\text{opt}}, \beta_{\text{opt}}\rangle$ and measure multiple times in the computational basis to extract the solution, which is the bit string with the highest probability. If there are n optimal solutions, they can be extracted by choosing the n bit strings with the highest probabilities.

A.2 Results

The implementation from the previous section is run on two different Quantum Inspire back-ends: the QX simulator back-end and the Starmon-5 QPU back-end. The cost function is optimized for 30 iterations using the constrained optimization by linear approximations (COBYLA) algorithm. For the quantum circuit executions during optimization 1024 shots were used, and for the final measurement 4096 shots were used. The results from the experiments are shown in Figure A.5. The QX simulator back-end manages to find a local minimum with a final approximation ratio of 0.76 and high probabilities of measuring the optimal solutions $z = 1010$ and $z = 0101$. The Starmon-5 QPU back-end reaches a final approximation ratio of 0.63 and manages to reach a final state which measures an optimal solution $z = 0101$ with high probability. The performance difference between the simulator and QPU back-end was expected: while HQCAs have some robustness against noise, their performance on real hardware will still be significantly worse from idealized noise-free simulations. Further performance improvement could be achieved on QPUs by choosing a classical optimizer that is more robust against noise as discussed in [73, 74], but such work is beyond the scope of this report.

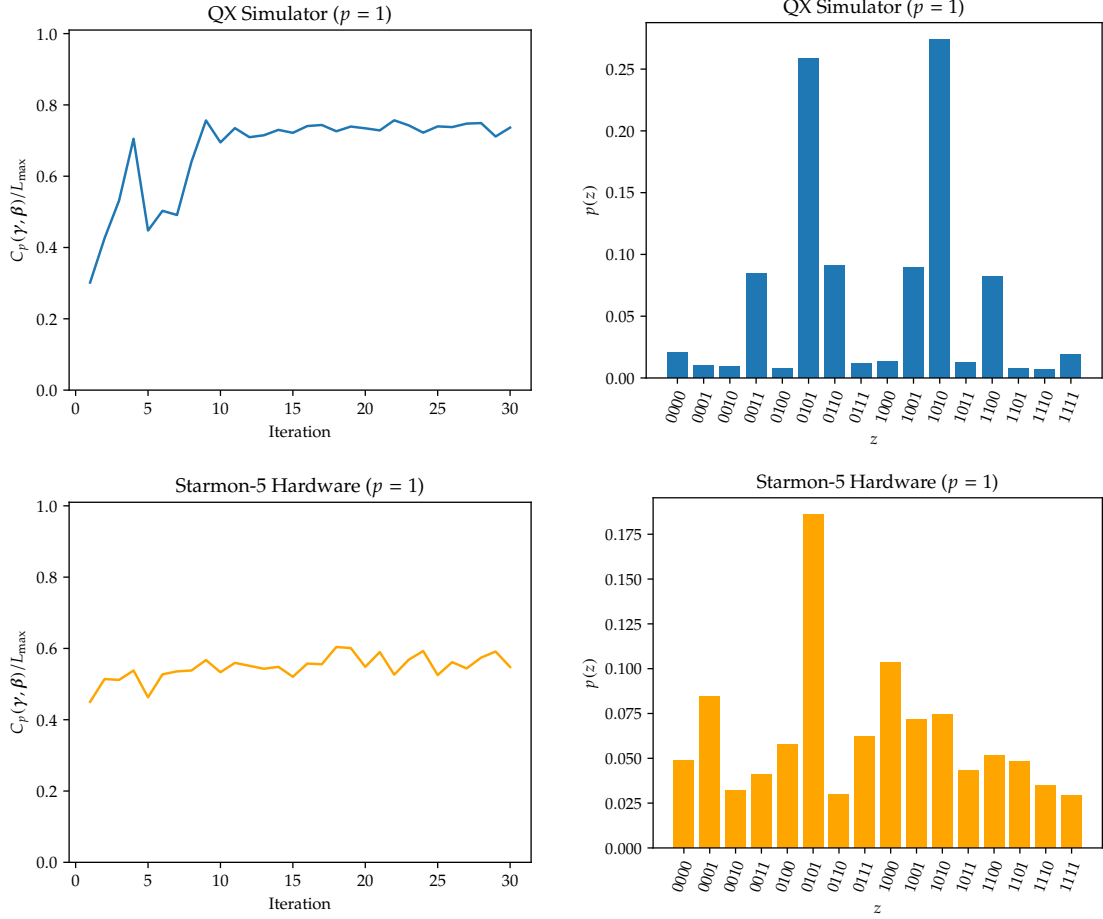


Figure A.5: Experimental results of the running QAOA to solve the Max-Cut problem on the graph from Figure A.2. The left column plots the approximation ratio $C_p(\gamma, \beta)/L_{\max}$ over 30 iterations, and the right column plots the final probabilities of measuring the possible solution bit strings after optimization. The top row contains results from the QX simulator back-end, and the bottom row contains the results from the Starmon-5 QPU back-end.