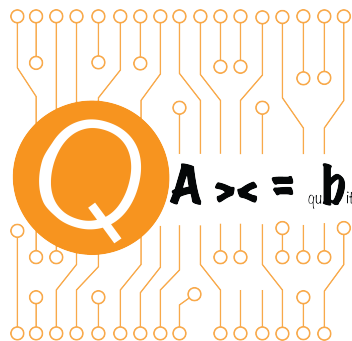


# Quantum Algorithms for Scientific Computation

Lin Lin and Nathan Wiebe

January 21, 2026



PRELIMINARY NOTES BEING CONTINUOUSLY UPDATED.



# Contents

<b>Part I. Background</b>	<b>5</b>
Chapter 1. Quantum advantage in scientific computation	7
1.1. Origin and Justification for Quantum Computing	7
1.2. Quantum speedup	10
1.3. Quantum advantage hierarchy	12
1.4. Quantum error correction and fault tolerant computation	15
1.5. Error accumulation mechanisms in classical and quantum computation	16
Chapter 2. Elements of quantum computation	21
2.1. Basic notation	21
2.2. Postulates of quantum mechanics	23
2.3. Density operator	34
2.4. Quantum circuit	38
2.5. Copy operation and no-cloning theorem	42
2.6. Deferred and implicit measurements	44
2.7. Sparse matrix, Majorana, fermionic, and bosonic operators	46
2.8. Selected Examples of Hamiltonians in Physics, Chemistry, and Optimization	49
<b>Part II. Foundation</b>	<b>57</b>
Chapter 3. Probability, quantum channel, and distances	59
3.1. Basic notions in probability theory	59
3.2. Quantum Channels	62
3.3. Distance between state vectors and unitaries	69
3.4. Distance between classical states and classical channels	72
3.5. Distance between quantum states	73
3.6. Distance between quantum channels	82
Notes and further reading	91
Chapter 4. Universality of quantum circuits	93
Chapter 5. Quantum processing of classical information	95
5.1. Reversible simulation of classical gates	95
5.2. Uncomputation	96
5.3. Fixed point number representation and quantum random access memory	99
5.4. Classical arithmetic operations	100
Notes and further reading	103

Chapter 6. Query complexity and quantum complexity theory	105
Chapter 7. Perturbation theory	107
Chapter 8. Statistical estimates	109
<b>Part III. Algorithm</b>	111
Chapter 9. Block encoding	113
9.1. Block encoding	113
9.2. Linear combination of unitaries	116
9.3. Block encodings of matrix additions and multiplications	120
9.4. Example: implementing generalized measurements	123
9.5. Example: Quantum error correction as block encoding	123
9.6. Query models for matrix entries	123
9.7. Block encoding of $s$ -sparse matrices	124
9.8. Hermitian block encoding	128
Notes and further reading	129
Bibliography	131
Index	135

## Part I

# Background

Part I of this book sets the stage for our exploration of quantum algorithms for scientific computation by asking two questions: why should we expect quantum computers to offer a computational advantage, and what are the basic mathematical and physical principles that govern them?

Chapter 1 tackles the first question. We begin by tracing the conceptual origins of quantum computing and formalizes the notion of quantum speedup. We then introduce a quantum advantage hierarchy, which classifies applications based on the strength of evidence for quantum speedup.

Chapter 2 addresses the second question by providing a concise overview of elements of quantum computation. We introduce the postulates of quantum mechanics, the circuit model, and the density operator formalism. We also cover concepts such as the no-cloning theorem and the principles of deferred and implicit measurement. The chapter concludes by introducing the operator formalisms for spin, fermionic, and bosonic systems, which are essential for describing the physical problems encountered in scientific applications, and presents several example Hamiltonians that will serve as recurring illustrations throughout the book.

## CHAPTER 1

# Quantum advantage in scientific computation

In this chapter, we trace the conceptual origins of quantum computing and explain how the physical nature of information suggests that quantum mechanics may offer computational power beyond classical Turing machines. We then formalize the notion of quantum speedup. Any claim of quantum advantage requires accounting for all relevant computational costs, including data input and output. To structure this assessment, we introduce a quantum advantage hierarchy that categorizes problems based on the existing evidence for significant speedups. The chapter concludes with a brief discussion of quantum error correction, and why exponentially large state spaces do not force exponential error accumulation: in fault-tolerant computation, it suffices to implement each gate to an accuracy that scales inversely with the gate count.

### 1.1. Origin and Justification for Quantum Computing

Our aim in this textbook is to provide a concrete understanding of not only how quantum algorithms work, but more importantly *why* they work and what impact scalable quantum computers are expected to yield in both the scientific and industrial worlds. Underlying this inquiry, however, is a deeper philosophical question about what it means to compute and why probing this question inevitably led to the idea of quantum computing.

Modern computer science traces its roots back to the early 20th century, with luminaries such as Alan Turing, John von Neumann, and Claude Shannon struggling to mathematically describe how information is stored and processed. Turing’s great realization was that all such computers could be mathematically modeled by an abstract device called a “Turing Machine”. The Turing machine was inspired strongly by the human “computers” (clerks) of the day: it possesses a tape for storing information and a read head that moves along the tape, updating the data on the tape in accordance with a stored program [Tur36].

John von Neumann is often credited with providing the first modern computer architecture that resembles modern computers, featuring dedicated memory, arithmetic and logic units, and input/output capabilities [VN93]. This architecture provided a far more realistic model of the postwar computers that were emerging, but conceptually these devices were no more powerful than the original Turing machine. Specifically, a machine is said to be “Turing Complete” if any function that a Turing machine can compute can be computed on the device. The von Neumann machine (given sufficient memory) can be shown to be Turing Complete, and in fact, a Turing machine can also simulate a machine implementing the von Neumann architecture. In this sense, the device is more than just Turing Complete: it is actually Turing Equivalent. Indeed, all known classical computational systems are Turing equivalent in this sense. This observation means that, effectively, every computational system in the universe could be understood as a Turing machine.

The formal study of algorithms revealed that not all tasks are fundamentally as easy for a Turing Machine. Some tasks, such as deciding whether a program halts, are strictly uncomputable [Tur36].

On the other hand, problems such as multiplying two  $n$ -bit numbers can be performed using a number of steps that scales polynomially in  $n$ . Still other problems, such as factoring an  $n$ -bit integer into a product of primes, can have their solution *verified* in polynomial time, but to date, no efficient algorithm has been found on a Turing machine that can *find* these factors in time that is polynomial in  $n$  (despite centuries of study). This suggested that a more fine-grained notion of computability needed to be considered than simply “computable” or “uncomputable”. Instead, it was seen to be useful to categorize computational tasks that can be computed on a Turing Machine using a polynomial number of operations as “efficiently computable” and all others as inefficient.

This categorization led to a bold hypothesis, which we will later criticize, known as the **Extended Church-Turing Thesis**. This statement says that any reasonable model of computing can be simulated using a polynomial number of computational steps by a probabilistic Turing machine. The example of von Neumann’s model of computing being simulatable in polynomial time by a Turing machine has indeed been reinforced by other models of computing based on physical phenomena, including billiard balls and the Game of Life. However, a challenge would emerge from an unlikely source: fundamental physics.

At the same time as computer science was being developed, a revolution was happening in physics. It had long been observed by physicists such as Planck and Einstein that classical physics could not be used to explain why heated objects (blackbodies) glowed red or how solar panels worked. Indeed, realistic models of these effects based on Newtonian principles failed to predict experimental observations. In the case of the stove elements, this failure was so radical that it predicted that infinite energy would be emitted by a stove burner (the “ultraviolet catastrophe”). A new type of model, formalized by von Neumann and others, was proposed to describe these systems that we now know as quantum mechanics (so named for its prediction that light should be emitted or absorbed in discrete quanta of energy). This language ultimately became the foundation of all fundamental physical law (gravitation being a notable exception).

Subsequent questions from Einstein, Podolsky, Rosen, and developments by Bell showed that quantum mechanics could not reasonably be described by classical local realism. Specifically, a phenomenon known as **entanglement**, which describes the correlations between measurement outcomes of coupled quantum systems, could not be described by classical mechanics without incorporating a non-local mechanism for updating measurement results. This work began to seriously question whether quantum systems could be plausibly described as mechanical systems. This, in turn, would much later be seen to question the Extended Church-Turing Thesis, as a Turing Machine is at its core a classical mechanical object that relies on local interactions.

A surprising feature of quantum mechanics is that its connection to computing seems to have taken several decades to be appreciated, despite us owing John von Neumann a great debt for formalizing both theories. With the benefit of hindsight, it is clear that with the appreciation of the fact that information is physical, quantum computing could have been developed as early as the 1940s.

The physical nature of information was elucidated most clearly by Shannon and Landauer. Shannon showed that the information content of a signal takes the same form as entropy, or disorder, in thermodynamics. Inspired by this connection, Shannon proposed that the two concepts were the same, establishing a link between his mathematical theory of information and thermal physics. Indeed, according to a widely circulated anecdote attributed to Shannon in an article by Tribus, von Neumann may have been agonizingly close to realizing the connection between physics and information processing [TM71]:



*“What’s in a name? In the case of Shannon’s measure the naming was not accidental. In 1961 one of us (Tribus) asked Shannon what he had thought about when he had finally confirmed his famous measure. Shannon replied: ‘My greatest concern was what to call it. I thought of calling it ‘information’, but the word was overly used, so I decided to call it ‘uncertainty’. When I discussed it with John von Neumann, he had a better idea. Von Neumann told me, ‘You should call it entropy, for two reasons. In the first place your uncertainty function has been used in statistical mechanics under that name. In the second place, and more importantly, no one knows what entropy really is, so in a debate you will always have the advantage.’”*

Indeed, Shannon’s work provided strong evidence that the two concepts are in fact the same and that thermodynamics had been telling us a secret lesson about information all along.

Landauer took this insight one step further by showing that thermodynamics places limitations on computers. Specifically, he showed that any computer that performs a calculation at finite temperature must pay an energy price for every bit of information erased to avoid violating the laws of thermodynamics [Lan61]. Similar work studying Maxwell’s Demon, a hypothetical agent that can raise and lower a gate that allows fast gas molecules through while blocking slow molecules, revealed that if the thermodynamic cost of measuring and computing were ignored, the laws of thermodynamics could be violated by such an agent [Ben87]. These works showed a strong link between information and physics and laid the foundation for the link to quantum computing that would soon follow.

It took the insight that information is physical to begin to motivate incorporating the formalism of quantum mechanics into the language of computer science. Quantum computing was born of this synthesis and was articulated independently by Manin [Man80] and Feynman [Fey82]. The justification that they had was the fact that the description of the state space of even small quantum systems scales exponentially with the number of quantum bits. This means that a naïve simulation of the laws of quantum mechanics would require exponentially more time on a classical computer than the physical system itself requires to evolve. This work opened the possibility that a computer that exploited the full capabilities of quantum mechanics may be, for certain problems, exponentially more powerful than the Turing machine. This in turn caused the scientific community to begin to doubt that the Extended Church-Turing Thesis holds, and now the belief that any realistic model of computing is polynomially equivalent to a quantum computer has become widespread after the discovery of the fast factoring algorithm of Shor [Sho99], the quantum simulation algorithms of Lloyd and others [Llo96], as well as the quantum advantage proposals of Aaronson and Arkhipov [AA11].

At a high level though, quantum computing suggests something potentially even stronger. If the Extended Church-Turing Thesis is replaced by a quantum version, then all of nature could be described or simulated in polynomial time by a massive quantum computer. In this sense, the strong link between information and physics reaches a crescendo with quantum computing, which suggests that all of physical law could be thought of as an algorithm that is run on a quantum computer, and the set of tasks that a quantum computer cannot perform efficiently are precisely those that nature also cannot solve at scale. For this reason, the search for exponential algorithmic advantage plays a central role in quantum computing, not only because it provides us with new opportunities for our computers, but also because it reveals the limitations that physical systems impose on information processing, and in turn, the limitations that information processing places on physical systems. Indeed, the main purpose of this text is to shed light on the origin and utility of quantum speedups for scientific applications.

### 1.2. Quantum speedup

The primary aim of exploring quantum computation is to attain a **quantum speedup** or **quantum advantage**, thereby enhancing problem-solving capabilities in scientific computation. At first glance, it seems that  $n$  qubits can be used to represent a superposition over  $2^n$  classical basis states, and significant quantum speedups should be expected everywhere. However, the situation is much more ambiguous: does the quantum algorithm require an exponential amount of classical information to be passed into the quantum computer? Does the quantum algorithm generate an exponential amount of information that needs to be extracted out of the quantum computer? If the size of the classical state space is  $2^n$ , is it mandatory for the classical algorithm to go through all states in order to find an approximate solution to a desired precision? If the size of the classical state space is only  $n$  but the computational cost of an existing algorithm is  $2^n$ , is it possible for a future classical algorithm to reduce this cost to  $\text{poly}(n)$ ? Readers may be curious about how to evaluate and answer these questions before dedicating substantial time to learning quantum computation. Indeed, these discussions can occur at a relatively broad level, largely circumventing the need for intricate quantum jargon.

One way to formulate the quantum speedup (as a function of the system size  $n$ ) is

$$(1.1) \quad \text{Quantum speedup} = \frac{\log(\min \text{Cost}(\text{classical}))}{\log \text{Cost}(\text{quantum})}.$$

The presence of the logarithm can be intuitively understood as follows. For a task with a “system size”  $n$ , assume that the classical and quantum costs are (asymptotically) proportional to  $n^{\alpha_c}$  and  $n^{\alpha_q}$ , respectively. Then as  $n \rightarrow \infty$ , the quantum speedup defined according to Eq. (1.1) is  $\alpha_c/\alpha_q$ . For instance, a *quadratic* quantum speedup means  $\alpha_c/\alpha_q = 2$ , a *cubic* quantum speedup means  $\alpha_c/\alpha_q = 3$ , and so on. If  $\alpha_c \rightarrow \infty$  as  $n \rightarrow \infty$  but  $\alpha_q$  remains bounded, the quantum speedup is *superpolynomial*. There is also a concept called “exponential quantum advantage” (EQA), which suggests that the classical cost increases at least exponentially in  $n$  but the quantum cost increases only polynomially.

Rigorous proof of EQA can be extraordinarily difficult for practical problems. For example, given two prime numbers  $p, q$ , the product  $m = p \cdot q$  can be easily carried out on a classical computer. However, if we are only given the integer  $m$ , finding the prime factors  $p, q$  can be very challenging. This is called the prime factorization problem and has wide applications in cryptography. The difficulty of the prime factorization problem can be measured in terms of the number of bits in  $m$ . An integer  $m$  can always be expressed in binary format. For instance,  $12 = 2^3 + 2^2$  can be represented as 1100 in binary format, where the number of bits  $n$  is 4. The most efficient classical algorithm, judged by asymptotic scaling in  $n$ , is the General Number Field Sieve method [Bri98]. The computational scaling is proportional to  $\exp[cn^{\frac{1}{3}}(\log n)^{\frac{2}{3}}]$ , which increases superpolynomially with  $n$ . Shor’s celebrated algorithm [Sho94, Sho99] addresses the same problem on a quantum computer, with its cost being proportional to  $n^2 \log n \log \log n$ , i.e., only polynomial in  $n$ . On one hand, this provides a very clean (and so far the cleanest) quantum solution with a significant quantum speedup that is superpolynomial in  $n$ . On the other hand, even for this problem, the speedup is not yet exponential in the strict sense above. For practical purposes, we will be (more than) content with a superpolynomial quantum speedup.

In principle, the classical cost should be minimized with respect to *all* classical algorithms, including algorithms that exist today, and those that will ever be developed in the future. A useful lower bound of the cost of classical algorithms may be obtained for some simple problems. However, this undertaking is exceedingly challenging for the majority of scientific computing problems. For

instance, we do not know whether the problem of prime factorization can or cannot be performed in polynomial time. Therefore, for practical purposes, we will further be satisfied with an estimate of  $\min \text{Cost}(\text{classical})$  by weighing both theoretical and empirical evidence, based on *existing* classical algorithms.

Although quantum mechanics is frequently described as a probabilistic theory, a key component is actually the quantum wavefunction (or quantum amplitude). This can be roughly equated to the square root of a probability density, along with phase information. This difference between probability density and quantum amplitude often forms the basis of the quadratic speedup, i.e.,  $\alpha_c/\alpha_q = 2$ . The most prominent example of this is Grover's algorithm for unstructured search (see ??). Although a quadratic speedup is valuable, it is unlikely that this speedup alone will be the most groundbreaking application of early fault-tolerant quantum computers. Hence, we use the loose term *significant* quantum speedup to refer to speedups greater than quadratic (such as cubic or quartic), or better, to superpolynomial speedups.

The quantum cost can be roughly calculated as the total gate complexity multiplied by the number of repetitions due to the measurement process. It is also conceptually useful to divide it into the following three components:

- (1) Input cost, or the cost for preparing the input quantum state. Without loss of generality, the quantum algorithm starts from a clean quantum state such as  $|0^n\rangle$ , and the input state to the quantum algorithm, denoted by  $|\psi_I\rangle$ , can be prepared using a unitary matrix  $U_I$  as  $|\psi_I\rangle = U_I |0^n\rangle$ . Then the input cost is the gate complexity for implementing  $U_I$ . Sometimes a quantum algorithm requires multiple accesses to the input oracle  $U_I$  in a coherent fashion. In this case, the input cost is given by the gate complexity for implementing  $U_I$  multiplied by the number of coherent initial state preparations.
- (2) Output cost, or the cost of quantum measurement. Without loss of generality, after an appropriate basis change the measurement can be taken to be performed on one or multiple qubits in the computational basis at the end of an algorithm. Then the output cost is the number of repetitions  $M$  needed to run the quantum algorithm.
- (3) Running cost, or the cost of coherently running the quantum algorithm once. This is given by the gate complexity for implementing the algorithm (excluding the cost for implementing  $U_I$ ).

One reason for separating the total gate complexity into the input cost and the running cost is that it allows us to distinguish the case when the overall cost is dominated by preparing the input, rather than by coherently executing the rest of the algorithm. In many settings, the input information is classical, and the nature of its complexity can be very different from that in the quantum algorithm. There is also an important scenario in which the input state  $|\psi_I\rangle$  is not generated by a known circuit  $U_I$ , but is produced by a quantum experiment. In this case, the relevant input cost is often the number of times the experiment must be repeated to prepare  $|\psi_I\rangle$  (a sample complexity), rather than the gate complexity of a circuit. For instance, **quantum learning theory** studies how efficiently one can infer properties of an unknown quantum state from state preparations and measurements. Throughout this book, we focus on computational tasks in which quantum and classical algorithms have access to the same amount of classical input information and are required to output classical information, and we will not discuss quantum learning theory in detail (except basic concepts such as parameter estimation in Chapter 8).

Ultimately, all quantum algorithms must output information that can be processed through classical means via quantum measurements. If the quantum state itself is the end product, the procedure to recover the quantum state on a classical computer is called quantum state tomography.

The cost of the state tomography procedure usually grows exponentially relative to the size of the quantum system. Therefore, it is unlikely that significant quantum speedup can be achieved for problems involving a tomography procedure on a large number of qubits. Instead, we should focus on problems whose end result can be obtained by measuring a small number of observables related to the quantum state to a desired accuracy, for which the measurement overhead can sometimes be reduced substantially.

In summary, a quantum computer should not be viewed as an all-purpose computational device destined to replace classical computers. Rather, it should be seen as an accelerator, capable of providing significant speedups for specific computational tasks. As emphasized in [Aar14], one must “read the fine print” when evaluating claims of quantum advantage. Several criteria must be met: the problem under consideration should be computationally intensive on classical hardware; the task must be solvable efficiently on a quantum device; and the overhead associated with data input and output (i.e., loading and extracting data) should not dominate the overall cost. Furthermore, several proposed quantum speedups for linear algebra and machine learning on classical data rely on strong data-access assumptions, and in some cases comparable scaling can be achieved by quantum-inspired classical algorithms under similar assumptions. Meeting all of these conditions is far from trivial. It represents a significant theoretical, experimental, and algorithmic challenge for the entire scientific community.

### 1.3. Quantum advantage hierarchy

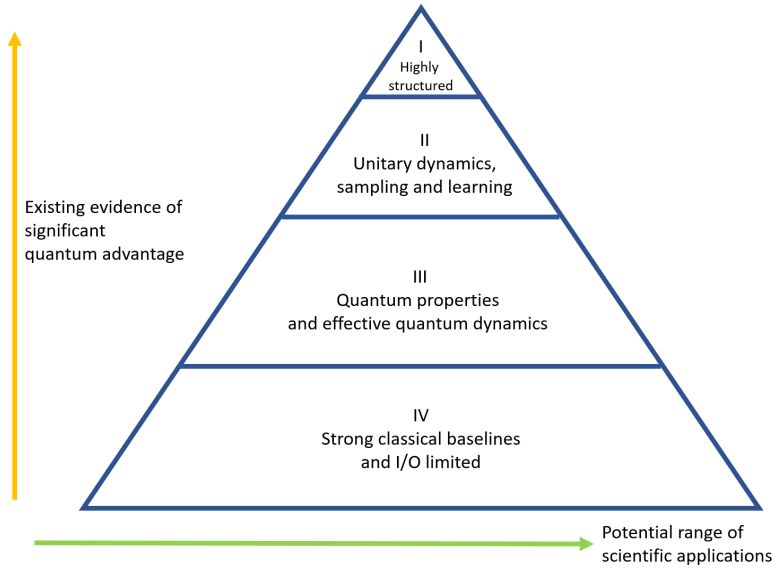


FIGURE 1.1. Quantum advantage hierarchy. The vertical level is determined by the most compelling application in each category, as demonstrated by the available evidence of quantum advantage.

Based on the aforementioned definition of quantum speedups, Fig. 1.1 organizes various quantum applications related to this book using a pyramid structure of 4 levels: Level I (Highly structured problems), Level II (Unitary dynamics, sampling, and learning problems), Level III (Quantum properties and effective dynamics problems), and Level IV (Strong classical baselines and I/O limited problems). Significant quantum speedups may exist across all levels. The vertical axis represents the *existing* amount of evidence supporting significant quantum speedups, while the horizontal axis represents the *potential* range of scientific applications. Besides gate complexity, for learning and sensing tasks the dominant cost can be the sample complexity (the number of experimental repetitions), which we treat as part of the running and output costs in this hierarchy. Now we give some examples at each level of the hierarchy, and the results are summarized in Table 1.1.

Level	Input Cost	Output Cost	Running Cost	Classical Cost	Examples
I Highly structured	✓	✓	✓	Provably expensive	Shor's algorithm for prime factorization and discrete logarithm, decoded quantum interferometry for structured optimization
II Unitary dynamics sampling and learning	✓	✓	✓	Empirically expensive	Hamiltonian simulation, random circuit sampling, learning with quantum memory
III Quantum properties and effective dynamics	?	?	✓	Empirically expensive	Ground state energy estimation, thermal state preparation, Green's function, open quantum system dynamics
IV Strong classical baselines and I/O limited	?	?	?	Efficient (except very large systems)	Classical partial differential equations, stochastic differential equations, unstructured optimization, classical machine learning

TABLE 1.1. Examples of problems in the quantum advantage hierarchy and existing amount of evidence justifying significant quantum speedups.

While prime factorization (and cryptography problems in general) are not typically classified as scientific computing problems, they occupy a unique position (Level I) at the peak of this hierarchy, and serves as a reference point for the ideal demonstration of quantum advantage in highly structured settings. These problems possess specific mathematical structures that allow quantum algorithms to bypass the exhaustive search often required by classical approaches. By describing the classical cost as “provably expensive,” we mean that the problem is hard under reasonable complexity-theoretic conjectures or relative to the best-known classical algorithms. For example, Shor's algorithm exploits the periodicity of the modular exponentiation function, a property related

to the hidden subgroup problem. Another reason for placing Shor’s algorithm at the top of the hierarchy is that these problems are “verifiable,” meaning a candidate solution can be efficiently checked by classical means (e.g., by multiplying the returned factors). The recently developed decoded quantum interferometry (DQI) algorithm [JSW<sup>+</sup>25] also solves highly structured optimization problems with superpolynomial speedups and has the potential to be classified in Level I. Unfortunately, such structures are rare in general scientific computing settings, as most problems in physics and engineering lack these clean, exploitable properties. To date, only a small number of applications have presented a comparable level of evidence supporting significant quantum speedups, and the list of credible candidates continues to evolve.

The most prominent example in Level II is the time evolution of a quantum state under a Hamiltonian, known as the Hamiltonian simulation problem. Many tasks in quantum physics and chemistry can be cast in this form. This category also includes sampling from unitary evolutions not explicitly defined by a Hamiltonian, such as random circuit sampling used in quantum supremacy experiments. For a physical Hamiltonian acting on  $n$  qubits, the description size is typically polynomial in  $n$ . We assume simple initial states, such as product states, which can be prepared with polynomial cost. The cost to simulate the dynamics for time  $t$  with precision  $\epsilon$  then scales as  $\text{poly}(n, t, 1/\epsilon)$ . Under these assumptions, no known classical algorithm is expected to reliably simulate generic many-body dynamics for long times. However, compared to Level I, the theoretical justification for speedup is often less rigorous, and verifying quantum advantage can be more difficult. For instance, verifying the output distribution of random circuit sampling typically demands exponential classical resources. Consequently, evidence for advantage relies heavily on the empirical hardness of classical simulation. Certain quantum learning tasks can demonstrate exponential advantage in sample complexity [HBC<sup>+</sup>22]. These advantages primarily stem from the quantum nature of the input data and the availability of quantum memory [CCHL22]. This differs from the computational tasks addressed in this book, which focus on problems with classical inputs and outputs. Furthermore, the exponential advantage assumes that we have zero knowledge about the quantum system being learned, which is often not the case in physical applications.

Level III of the hierarchy includes a large class of problems in quantum physics, quantum chemistry, and materials science. By “quantum properties,” we refer to static characteristics of the system, such as ground state energy, excited state energy, stationary states, and spectral properties. By “effective dynamics,” we refer to processes that are not natively unitary, such as open system dynamics involving dissipation or imaginary time evolution used for cooling. Compared to Level II problems, the mapping from these non-unitary objects to the unitary logic of quantum hardware is indirect. This mapping often introduces overheads, such as the need for linear combinations of unitaries, post-selection, or many ancillary qubits. The amount of information that needs to be extracted from the quantum computer can be comparable to that in the quantum dynamics simulation and is at most polynomial in  $n$ . In the case of the ground-state energy estimation, the situation is even clearer since we only need to estimate a single number as the output. Compared to unitary dynamics, there exist a much larger number of powerful classical algorithms for these tasks. These are approximate methods and often cannot be used to converge to the true solution to arbitrary precision. However, for many practical problems, they have been shown to be sufficiently accurate. Input cost is also a major factor placing these problems at Level III. For example, ground state estimation often requires a good initial guess (an ansatz) to succeed; generating this ansatz can be computationally expensive or physically difficult, sometimes leading to QMA-hard bottlenecks. Finally, we may design quantum algorithms to solve problems that are entirely classical. For instance, we can consider quantum solvers for classical partial differential equations

(PDEs), stochastic differential equations, unstructured optimization problems, and sampling tasks. These cover a large variety of problems in scientific computing. However, many such classical problems fall into Level IV because they have strong classical baselines and/or are I/O limited. For example, many PDEs on a grid of size  $N$  can be solved classically in time polynomial in  $N$  (often approximately linear in  $N$  using fast algorithms). Even if a quantum algorithm offers a speedup in the processing stage, it faces the “I/O limit”: merely loading an arbitrary input vector of size  $N$  into the quantum state takes time linear in  $N$ , which negates any potential exponential speedup. One exception arises when the classical data possesses significant structure that allows for efficient loading; a potential advantage in this regime was recently demonstrated for a quantum solver of a large number of classical oscillators [BBK<sup>+</sup>23]. Regarding unstructured search problems, while Grover’s algorithm provides a quadratic speedup, this is often insufficient to overcome the significant constant-factor overheads of fault-tolerant quantum error correction compared to highly optimized classical heuristics. Thus, while the range of applications is vast, securing an end-to-end advantage is difficult. That being said, many cryptography problems can be formulated as classical optimization problems, and the next breakthrough in quantum algorithms *may* emerge from classical problems again.

The ongoing evaluation and pursuit of quantum advantage is a rapidly developing field. When discussing applications, we will only scratch the surface of the potential indications of quantum advantage by examining aspects such as quantum input cost, output cost, running cost, and the cost of classical algorithms, wherever possible. This approach is intended to encourage readers to seek out these elements in their own research. However, it is important to understand that the findings presented, while based on existing literature, are far from exhaustive or conclusive. The rapid pace of advancements means that future developments could significantly alter the current understanding and conclusions.

#### 1.4. Quantum error correction and fault tolerant computation

All previous discussions assume that quantum operations can be perfectly performed. To this end, quantum error correction is necessary. The threshold theorem [ABO97] is a central result in the field of quantum error correction. The theorem essentially states that if the error rate of quantum operations (including gates and measurements) is below a certain threshold value (around 0.001, though the precise value depends on the detailed assumptions), then it is possible to perform quantum computation for an arbitrary length of time with arbitrarily high accuracy (see [NC00, Section 10.6]).

**THEOREM 1.1 (Threshold theorem).** *There exists an error threshold  $p_t > 0$ . If the physical error rate  $p$  per gate operation satisfies  $p < p_t$ , there exists a quantum error correction scheme such that the logical error rate  $q$  can be made as small as desired. In other words,  $q = \mathcal{O}((p/p_t)^\ell)$  for any positive integer  $\ell$ .*

We will not study the details of quantum error correction in this book. In classical computing, modern algorithm design generally does not take error correction into account. Similarly, in the long term, quantum error correction is expected to be largely a separate issue from the design of quantum algorithms. We always assume quantum error correction protocols have been implemented, physical noise has been eliminated, and the resulting quantum computer is **fault-tolerant**. For the purpose of this book, all errors come from either *approximation errors* at the mathematical level, or *Monte Carlo errors* in the readout process due to the probabilistic nature of the measurement process.

Quantum error correction is a dynamic and rapidly progressing field, and will significantly impact the development and potential of quantum algorithms, and the landscape of quantum computing. On a very coarse scale, we can categorize quantum algorithms based on the type of quantum computer architecture they are designed for.

- (1) Noisy intermediate-scale quantum (NISQ) computers: These devices represent the current state of quantum computing technology. Characterized by a relatively small number (tens to a few hundreds) of physical qubits, these systems are prone to errors and lack full error correction capabilities. Quantum algorithms designed for NISQ devices, such as the Variational Quantum Eigensolver (VQE), need to be error resilient and must be capable of delivering meaningful results despite the presence of noise. Most of this book will not discuss NISQ algorithms.
- (2) Fully fault-tolerant quantum computers: These are the ideal, long-term goal of quantum computing research. In these systems, quantum error correction protocols are fully implemented, allowing quantum algorithms to run for long durations without being overwhelmed by errors. This architecture will enable the execution of complex algorithms that require a large number of qubits and gate operations. Many of the algorithms discussed in this book are designed for this type of architecture. At the current stage, the goal of many fully fault-tolerant quantum algorithms is to minimize the total cost (in an *asymptotic* sense with respect to certain parameters, such as precision, system size etc.) for solving a given task.
- (3) Early fault-tolerant quantum computers: This category represents a transitional phase between NISQ devices and fully fault-tolerant quantum computers. These systems would implement some form of quantum error correction, but they may have constraints such as a very limited number of logical ancilla qubits. This means that they can only run quantum algorithms within a certain complexity limit. Despite these constraints, early fault-tolerant quantum computers provide an opportunity to test and refine fault-tolerant designs and protocols, and to run quantum algorithms that are beyond the reach of NISQ devices but do not require the full capabilities of fault-tolerant quantum computers. Some of the algorithms in this book take such constraints into account and can be suitable on early fault-tolerant quantum computers.

### 1.5. Error accumulation mechanisms in classical and quantum computation

Quantum computation aims at processing objects whose natural dimension is exponential, such as vectors in  $\mathbb{C}^{2^n}$  and matrices of size  $2^n \times 2^n$ . No computation can be carried out exactly, so will the error also accumulate exponentially with the system size? If that were the case, then quantum algorithms would become useless precisely in the regime where they are designed to operate. In this section we give a bird's-eye view of the relevant error accumulation mechanisms.

At first glance, deterministic numerical computation can look discouraging in this respect. Even a basic task such as forming an inner product involves many elementary operations, and Example 1.2 shows a worst-case bound for the accumulated rounding effects that is proportional to  $N = 2^n$ .

However, scientific computation has long dealt with exponentially large state spaces without requiring errors to grow linearly in the dimension. Randomized algorithms on  $n$  bits evolve a probability distribution on a space of size  $N = 2^n$ , yet the accuracy of the computation is governed by how many transition steps are composed, not by  $N$  itself: if each step is implemented to accuracy  $\epsilon$ , then the overall error is at most  $K\epsilon$ , where  $K$  is the number of steps (see Proposition 3.29).



Quantum computation behaves in the same way at the level of circuit synthesis: a quantum algorithm is a product of elementary unitaries, and the accumulated implementation error is controlled by the number of gates. In particular, if the gate count is  $K$  and each gate can be implemented to precision  $\epsilon/K$ , then the final error is  $\mathcal{O}(\epsilon)$ . In the fault-tolerant setting assumed above, achieving such per-gate accuracy is a realistic requirement, and the overhead of approximating elementary unitaries to a desired precision is discussed later (see Chapter 4). The distance notions used to make these comparisons precise are developed in Chapter 3, and the Monte Carlo errors arising at readout are discussed further in Chapter 8.

**1.5.1. Deterministic classical computation.** Modern scientific computation on classical computers is based on floating point arithmetic operations, which express a number in scientific notation. For instance, the number  $-0.271828 \times 10^5$  involves a sign ( $-$ ), fraction (271828), base (10), and exponent (5). In binary floating point, one stores a sign bit together with a fixed-length exponent and fraction. For instance, the IEEE single precision uses 1 bit for the sign, 8 bits for the exponent, and 23 bits for the fraction (32 bits long). The IEEE double precision uses 1 bit for the sign, 11 bits for the exponent, and 52 bits for the fraction (64 bits long). For instance, a double precision ranges from  $2^{-1022}$  to  $2^{1023}$ , or about  $10^{-308}$  to  $10^{308}$ . Numbers outside this range yield underflow or overflow error and need to be handled separately. This is much more efficient than the fixed point number representation (see Section 5.3), which would require more than 2046 bits (i.e., more than 2046 logical qubits for a single number) to cover the same range of numbers.

The basic assumption is that any real number  $a$  should be represented by  $\text{fl}(a)$  using a given number of bits. Similarly, any binary operation  $a \odot b$  should be represented by  $\text{fl}(a \odot b)$ , where  $\odot$  is one of the four elementary binary operations  $+$ ,  $-$ ,  $*$ ,  $/$ . The difference  $a \odot b - \text{fl}(a \odot b)$  is called the roundoff error. When the number is rounded correctly, i.e.,  $\text{fl}(a \odot b)$  is a nearest floating point number to  $a \odot b$ , we have

$$(1.2) \quad \text{fl}(a \odot b) = (a \odot b)(1 + \delta),$$

where  $|\delta|$  is upper bounded by  $\epsilon_{\text{mach}}$  (called the machine precision).

**Example 1.2.** Given  $u, v \in \mathbb{R}^N$ , consider the error accumulation of computing an inner product  $\sum_{i=1}^N u_i v_i$ . The error from each operation in the floating-point arithmetic needs to be counted separately. The floating-point representation of a product  $u_i v_i$  is given by  $u_i v_i(1 + \epsilon_i)$ , where  $|\epsilon_i| \leq \epsilon_{\text{mach}}$ , and  $\epsilon_{\text{mach}}$  is the machine epsilon.

However, when summing these products, there is an additional error introduced at each addition step. Let us denote by  $\delta'_j$  the relative rounding error incurred when adding the  $j$ -th term (so  $|\delta'_j| \leq \epsilon_{\text{mach}}$ ). Then the partial sums satisfy

$$(1.3) \quad \text{fl}(s_{j-1} + u_j v_j(1 + \epsilon_j)) = (s_{j-1} + u_j v_j(1 + \epsilon_j))(1 + \delta'_j),$$

where  $s_{j-1}$  denotes the computed partial sum from the previous step. After summing over all  $N$  terms, we may write

$$(1.4) \quad 1 + \delta_i := (1 + \epsilon_i) \prod_{j=i+1}^N (1 + \delta'_j).$$

Therefore if overflow or underflow does not occur, then

$$(1.5) \quad \text{fl}\left(\sum_{i=1}^N u_i v_i\right) = \sum_{i=1}^N u_i v_i(1 + \delta_i), \quad |\delta_i| \leq (1 + \epsilon_{\text{mach}})^N - 1 \leq e^{N\epsilon_{\text{mach}}} - 1.$$

◇

When  $N\epsilon_{\text{mach}} < 1$ , we have  $|\delta_i| \leq 2N\epsilon_{\text{mach}}$ . So the error grows linearly in  $N$ . This is due to the step of adding  $N$  numbers following a linear order. For computing the inner product, the error accumulation in the summation step can be significantly reduced using a technique called the pair summation (or cascade summation) to  $\mathcal{O}((\log N)\epsilon_{\text{mach}})$ . However, such a more accurate summation method is more difficult to implement in broader scenarios such as matrix-matrix multiplication. For most of the tasks, the  $\text{poly}(N)$  factor in the error accumulation is unavoidable. For instance, for solving a triangular linear system, the error accumulation is  $\mathcal{O}(N\epsilon_{\text{mach}})$  [GVL13, Chapter 3.1]. For Gaussian elimination (or  $LU$  factorization), standard backward-error bounds involve the growth factor  $\rho$  and scale polynomially in  $N$ , typically of order  $\mathcal{O}(N\rho\epsilon_{\text{mach}})$  [GVL13, Chapter 3.4].

That being said, not all deterministic computations involving vectors in  $\mathbb{C}^N$  necessarily exhibit a  $\text{poly}(N)$  accumulation of numerical error. Error accumulation is governed not by the ambient dimension  $N$  itself, but by the number of elementary operations performed. For instance, tensor network methods provide settings in which certain computations on structured vectors in  $\mathbb{C}^N$  can be carried out using only  $\text{poly}(n)$  operations, where  $N = 2^n$ . We will not discuss tensor network methods in this book, and classical probabilistic computation provides a more direct analogy to quantum computation for tackling high dimensional problems, as discussed next.

**1.5.2. Probabilistic classical computation.** A probabilistic computation on  $n$  bits evolves a probability distribution on a space of size  $N = 2^n$ , and hence it can be described by a vector in  $\mathbb{R}^N$  acted on by stochastic matrices. The ambient dimension is exponential in  $n$ , but the computation is specified by a sequence of local update rules. As a result, neither the cost nor the accumulated implementation error needs to scale exponentially in  $N$ . This viewpoint also extends to the comparison between quantum and classical algorithms: a probability distribution can be viewed as a special quantum state, and a transition matrix can be associated with a special quantum channel (see Section 3.2).

If we can implement each transition matrix to precision  $\epsilon$ , the global error of the overall transition matrix grows at most linearly with respect to the number of transition matrices and is at most  $1, K\epsilon$  (see Proposition 3.29). Equivalently, if the gate complexity is  $K$  and we can implement each transition matrix to precision  $\epsilon/K$ , then the final error is upper bounded by  $\epsilon$ , independent of  $N$ . Compared to deterministic classical algorithms, randomized algorithms introduce another error mechanism: even when the transition rule is specified, one often estimates quantities of interest by sampling, and the output is therefore subject to Monte Carlo fluctuations. For example, estimating an expectation value by  $N_s$  independent samples typically incurs an error of order  $\mathcal{O}(N_s^{-1/2})$ , independent of the size of the underlying sample space. The statistical side of this issue is discussed further in Chapter 8.

**1.5.3. Quantum computation.** Quantum algorithms are designed to handle objects of size  $N = 2^n$  without explicitly storing  $N$  numbers. As in probabilistic computation, error accumulation depends on how many steps are composed and on the metric used to compare channels (see Chapter 3), and they do not introduce an explicit dependence on  $N$ .

Every quantum circuit can be represented by a unitary  $U$ , decomposed into a series of simpler unitaries as  $U = U_K \cdots U_1$ . Each  $U_i$  can only be implemented approximately by some  $\tilde{U}_i$  to precision  $\epsilon$ . The implementation cost of each simple unitary is independent of the Hilbert space dimension  $N$  (see Chapter 4). This implies that for any vector  $|\psi\rangle$  of size  $N$ , the error between  $U_i|\psi\rangle$  and  $\tilde{U}_i|\psi\rangle$  is less than  $\epsilon$  with no explicit dependence on  $N$ .

If we can implement each local unitary to precision  $\epsilon$ , the global error grows at most *linearly* with respect to the number of gates and is at most  $K\epsilon$  (see Proposition 3.21). In other words, if the gate complexity is  $K$  and we can implement each gate to precision  $\epsilon/K$ , then the final error is upper bounded by  $\epsilon$  and is independent of  $N$ . The same statement holds for quantum channels (see Section 3.6).



## CHAPTER 2

# Elements of quantum computation

This chapter lays the groundwork for our journey into quantum algorithms for scientific computation. We will review the mathematical and physical principles that underpin quantum computing. While we assume a basic familiarity with quantum mechanics, our focus will be on establishing the specific concepts and notational conventions used throughout this book. This chapter is not intended as a comprehensive introduction to quantum computing, but rather as a targeted primer on the tools we will need to build and analyze sophisticated quantum algorithms. For a more comprehensive introduction to quantum computation, we refer the reader to standard textbooks such as [NC00, Wat18].

We start with the postulates of quantum mechanics, introducing the Dirac notation and the core principles governing quantum states and their evolution. We then move to the language of quantum circuits, which greatly simplifies the tensor manipulations inherent in multi-qubit systems. To handle scenarios involving noise and subsystems, we introduce the density operator formalism. We will also discuss the no-cloning theorem, which forbids the copying of arbitrary quantum states, and the principles of deferred and implicit measurement, which offer flexibility in circuit design. The latter part of the chapter introduces the representation of structured matrices, including sparse matrices and operators from fermionic and bosonic systems. We conclude with a selected list of Hamiltonians from physics, chemistry, and optimization that will serve as motivating examples in our exploration of quantum simulation and other applications.

### 2.1. Basic notation

The sets of real and complex numbers are denoted by  $\mathbb{R}$  and  $\mathbb{C}$ , respectively. For a complex number  $c \in \mathbb{C}$ , the notation  $\bar{c}$  or  $c^*$  denotes its complex conjugate.

A complex vector  $v$  of size  $N$  is an  $N$ -tuple of complex numbers, written as  $v \in \mathbb{C}^N$ , with its  $j$ -th component denoted by  $v_j$ . By default, we use 0-based indexing, that is,  $j \in [N] := \{0, \dots, N-1\}$ . When 1-based indexing is used, we will explicitly write  $j = 1, \dots, N$ .

The **vector 2-norm** of  $v$  is denoted by  $\|v\| = \sqrt{\sum_{i \in [N]} |v_i|^2}$ . Unless otherwise specified, a vector  $v \in \mathbb{C}^N$  is considered unnormalized. A nonzero, normalized vector (viewed as a pure quantum state) is written as  $|v\rangle = v/\|v\|$ . To emphasize that a vector is unnormalized, we sometimes use the notation  $|v\rangle_\times$ .

A matrix  $A$  of size  $M \times N$  is denoted by  $A \in \mathbb{C}^{M \times N}$ , and its  $(i, j)$ -th entry is  $A_{ij}$  or  $a_{ij}$ . For  $A \in \mathbb{C}^{M \times N}$ , the complex conjugate of  $A$ , denoted by  $\bar{A}$  or  $A^*$ , is obtained by replacing each entry of  $A$  with its complex conjugate. The inverse of  $A$  (if  $A$  is invertible) is denoted by  $A^{-1}$ . The transpose of  $A$  is denoted by  $A^\top$ . The Hermitian conjugate (or adjoint) of  $A$ , denoted by  $A^\dagger$ , is the complex conjugate of the transpose of  $A$ , which can be expressed as  $A^\dagger = (A^\top)^*$ . A matrix  $A$  is **Hermitian** if it is equal to its Hermitian conjugate, i.e.,  $A = A^\dagger$ . A matrix  $A$  is **normal** if it commutes with its Hermitian conjugate, i.e.,  $AA^\dagger = A^\dagger A$ . A matrix  $U$  is unitary if its Hermitian

conjugate is its inverse, i.e.,  $U^\dagger = U^{-1}$ . The set of all  $N \times N$  unitary matrices forms the unitary group, denoted by  $U(N)$ . The set of all  $N \times N$  unitary matrices with determinant 1 forms the special unitary group, denoted by  $SU(N)$ .

If all eigenvalues of a Hermitian matrix  $A \in \mathbb{C}^{N \times N}$  are nonnegative,  $A$  is called a **positive semidefinite** matrix, or **positive operator**, denoted by  $A \succeq 0$ . The notation  $A \succeq B$  means  $A - B \succeq 0$ , and  $A \preceq B$  means  $B \succeq A$ . Similarly, if all eigenvalues of  $A$  are positive, then  $A$  is called a **positive definite** matrix, denoted by  $A \succ 0$ . The notation  $A \succ B$  means  $A - B \succ 0$ .

The **operator norm** (also called **induced vector 2-norm**)<sup>1</sup> of a matrix  $A$  is

$$(2.1) \quad \|A\| := \sup_{\|v\|=1} \|Av\|.$$

In quantum information theory, it is useful to consider the **Schatten  $p$ -norm** of  $A$ :

$$(2.2) \quad \|A\|_p := \left( \text{Tr}(A^\dagger A)^{\frac{p}{2}} \right)^{\frac{1}{p}}, \quad p \geq 1.$$

The particularly useful one is the **Schatten 1-norm** (also called the **trace norm**)

$$(2.3) \quad \|A\|_1 := \text{Tr} \sqrt{A^\dagger A}.$$

For instance, any quantum state (density operator)  $\rho$  is normalized with respect to the trace norm, i.e.,  $\|\rho\|_1 = 1$ . Furthermore, the Schatten  $\infty$ -norm  $\|A\|_\infty$  can be shown to coincide with the operator norm  $\|A\|$ . Many readers may not be familiar with the Schatten norms. We will discuss these norms in detail in Chapter 3.

We adopt the following **asymptotic notations**: Let  $\mathbb{R}_+$  be the set of positive real numbers. Consider two functions  $f : \mathbb{R} \rightarrow \mathbb{C}$  and  $g : \mathbb{R} \rightarrow \mathbb{R}_+$ . For any  $a \in \mathbb{R} \cup \{\pm\infty\}$ , if  $\limsup_{x \rightarrow a} \frac{|f(x)|}{g(x)} < \infty$ , then we write  $f(x) = \mathcal{O}(g(x))$  as  $x \rightarrow a$ , or simply  $f = \mathcal{O}(g)$ <sup>2</sup> when  $x \rightarrow a$  is clear from the context. We write  $f = \Omega(g)$  if  $g = \mathcal{O}(f)$ ;  $f = \Theta(g)$  if  $f = \mathcal{O}(g)$  and  $g = \mathcal{O}(f)$ . Note that  $\mathcal{O}(g)$  can also be interpreted as a set, so it is also valid to write  $f \in \mathcal{O}(g)$ . Similarly we may write  $f \in \Omega(g)$ ,  $f \in \Theta(g)$  etc.

The notation  $\tilde{\mathcal{O}}, \tilde{\Omega}, \tilde{\Theta}$  are used to suppress subdominant polylogarithmic factors. Specifically,  $f = \tilde{\mathcal{O}}(g)$  if  $f = \mathcal{O}(g \text{ polylog}(g))$ ;  $f = \tilde{\Omega}(g)$  if  $f = \Omega(g \text{ polylog}(g))$ ;  $f = \tilde{\Theta}(g)$  if  $f = \Theta(g \text{ polylog}(g))$ . Note that these tilde notations usually do not suppress dominant polylogarithmic factors. For instance, if  $f = \mathcal{O}(\log g \log \log g)$ , then we write  $f = \tilde{\mathcal{O}}(\log g)$  instead of  $f = \tilde{\mathcal{O}}(1)$ . However, for simplicity of presentation, we may sometimes use the notation  $\tilde{\mathcal{O}}$  more casually to suppress dominant polylogarithmic factors. When we do so, we will make an explicit mention of this usage.

Throughout the book, the natural logarithm is denoted by  $\ln$ , and is sometimes written as  $\log$  without an explicit base when the context is clear. The logarithm to base 2 is denoted by  $\log_2$ . When  $N$  denotes the dimension of  $\mathbb{C}^N$ , and the notations  $N$  and  $n$  appear together, it is usually assumed that  $N = 2^n$  for some positive integer  $n$ , referred to as the number of quantum bits (or **qubits**). Additional notations will be introduced in the book as needed.

<sup>1</sup>In matrix analysis, the operator norm is sometimes denoted by  $\|A\|_2$  to indicate that this is the induced vector 2-norm. More generally, the induced vector  $p$ -norm is  $\|A\|_p = \sup_{\|x\|_p=1} \|Ax\|_p$  where  $\|x\|_p = (\sum_i |x_i|^p)^{1/p}$ . For example, the induced vector 1-norm is  $\|A\|_1 = \sup_{\|x\|_1=1} \|Ax\|_1 = \max_j \sum_i |a_{ij}|$ . This book **does not** adopt such a notation.

<sup>2</sup>Sometimes  $\mathcal{O}(g)$  is treated as a set of functions, and by this interpretation we can equivalently write  $f \in \mathcal{O}(g)$ .

## 2.2. Postulates of quantum mechanics

This section encapsulates some of the most important postulates of quantum mechanics. All postulates concern finite dimensional, closed quantum systems (i.e., systems isolated from environments). For more details, we refer readers to [NC00, Section 2.2].

### 2.2.1. State space postulate.

**Definition 2.1** (Hilbert space). *A (complex) Hilbert space denoted by  $\mathcal{H}$  is a complex vector space equipped with an inner product  $\langle \cdot | \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{C}$  that satisfies the following properties for all  $x, y, z \in \mathcal{H}$  and all  $\alpha, \beta \in \mathbb{C}$ :*

- (1) (Conjugate Symmetry)  $\langle x | y \rangle = \overline{\langle y | x \rangle}$ .
- (2) (Linearity in the second argument)  $\langle z | \alpha x + \beta y \rangle = \alpha \langle z | x \rangle + \beta \langle z | y \rangle$ .
- (3) (Positive-definiteness)  $\langle x | x \rangle \geq 0$  with equality if and only if  $x = 0$ .

Furthermore,  $\mathcal{H}$  is complete with respect to the norm induced by the inner product, where the norm of a vector  $x \in \mathcal{H}$  is given by  $\|x\| = \sqrt{\langle x | x \rangle}$ .

The state space postulate assumes that the set of all quantum states of a quantum system, called the state space, is a Hilbert space. If the state space  $\mathcal{H}$  is finite dimensional, it is isomorphic (i.e., there is a one-to-one mapping) to some  $\mathbb{C}^N$ , written as  $\mathcal{H} \cong \mathbb{C}^N$ . Throughout the book, unless otherwise specified, we only consider finite dimensional Hilbert spaces. A **state vector** (also called ket vector, wavefunction, or pure quantum state)  $|\psi\rangle \in \mathcal{H}$  can be identified with a column vector in  $\mathbb{C}^N$

$$(2.4) \quad \psi = \begin{pmatrix} \psi_0 \\ \psi_1 \\ \vdots \\ \psi_{N-1} \end{pmatrix}.$$

Let  $\{e_i\}$  be the standard basis of  $\mathbb{C}^N$ . The  $i$ -th entry of  $\psi$  can be written as an inner product  $\psi_i = \langle e_i | \psi \rangle$ . We also use the Dirac notation, which uses  $|\psi\rangle$  to denote a quantum state. We further postulate that two state vectors  $|\psi\rangle$  and  $c|\psi\rangle$  for some  $0 \neq c \in \mathbb{C}$  always refer to the same physical state. Hence without loss of generality we always assume  $|\psi\rangle$  is normalized to be a unit vector, i.e.,  $\langle \psi | \psi \rangle = 1$ . Restricting to normalized state vectors, the complex number  $c = e^{i\theta}$  for some  $\theta \in [0, 2\pi)$  is called the global phase factor.

Throughout the book, unless otherwise specified, an unnormalized state vector is often denoted by  $\psi$  without the ket notation  $|\cdot\rangle$ , and  $|\psi\rangle := \psi / \|\psi\|$  denotes the normalized counterpart.

The bra vector  $\langle \psi |$  can be interpreted as a linear functional on  $\mathcal{H}$ , which maps any  $|\varphi\rangle \in \mathcal{H}$  to a complex number  $\langle \psi | \varphi \rangle$ . When  $\mathcal{H} = \mathbb{C}^N$ , we have  $\langle \psi | \varphi \rangle = \sum_{i \in [N]} \overline{\psi_i} \varphi_i$ . It can be identified with a row vector, which is the Hermitian conjugate of the column vector  $\psi$ :

$$(2.5) \quad \psi^\dagger = (\overline{\psi_0} \quad \overline{\psi_1} \quad \cdots \quad \overline{\psi_{N-1}}).$$

The set of all bra vectors, or linear functionals on  $\mathcal{H}$ , is denoted by  $\mathcal{H}^*$ <sup>3</sup>.

Given a state space  $\mathcal{H}$ , let  $L(\mathcal{H})$  denote the set of all linear operators on  $\mathcal{H}$ . When  $\mathcal{H} = \mathbb{C}^N$ ,  $L(\mathbb{C}^N)$  can be identified with the set of  $N \times N$  matrices, denoted by  $\mathbb{C}^{N \times N}$ . The ketbra notation  $|\psi\rangle\langle\varphi|$  is an element in  $L(\mathcal{H})$ , which maps any vector  $|\xi\rangle \in \mathcal{H}$  to another state vector in  $\mathcal{H}$  as

<sup>3</sup>The star  $\star$  acting on a vector space does not mean the complex conjugation of  $\mathcal{H}$ . This notation is only used occasionally in the book. A Hilbert space satisfies  $\mathcal{H} \cong \mathcal{H}^*$  by the Riesz representation theorem.

$|\psi\rangle\langle\varphi|\xi\rangle$ . The matrix representation of  $|\psi\rangle\langle\varphi|$  is the product of the column vector  $\psi$  and the row vector  $\varphi^\dagger$ , i.e.,  $\psi\varphi^\dagger \in \mathbb{C}^{N \times N}$ .

**Example 2.2** (Single qubit system and Bloch sphere). A (single) qubit corresponds to a state space  $\mathbb{C}^2$ . We also define

$$(2.6) \quad |0\rangle = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad |1\rangle = \begin{pmatrix} 0 \\ 1 \end{pmatrix}.$$

Since the state space of the spin- $\frac{1}{2}$  system is also isomorphic to  $\mathbb{C}^2$ , this is also called the single spin system, where  $|0\rangle, |1\rangle$  are referred to as the spin-up and spin-down state, respectively. A general state vector in  $\mathbb{C}^2$  takes the form

$$(2.7) \quad |\psi\rangle = a|0\rangle + b|1\rangle = \begin{pmatrix} a \\ b \end{pmatrix}, \quad a, b \in \mathbb{C},$$

and the normalization condition implies  $|a|^2 + |b|^2 = 1$ . So we may rewrite  $|\psi\rangle$  as

$$(2.8) \quad |\psi\rangle = e^{i\gamma} \left( \cos \frac{\theta}{2} |0\rangle + e^{i\varphi} \sin \frac{\theta}{2} |1\rangle \right), \quad \theta, \varphi, \gamma \in \mathbb{R}.$$

If we ignore the irrelevant global phase  $e^{i\gamma}$  (which also absorbs a minus sign in the coefficient of  $|0\rangle$ ), then it holds

$$(2.9) \quad |\psi\rangle = \cos \frac{\theta}{2} |0\rangle + e^{i\varphi} \sin \frac{\theta}{2} |1\rangle, \quad 0 \leq \theta \leq \pi, 0 \leq \varphi < 2\pi.$$

So we may identify each single qubit quantum state with a unique point on the unit three-dimensional sphere (called the **Bloch sphere**) as

$$(2.10) \quad \mathbf{a} = (\sin \theta \cos \varphi, \sin \theta \sin \varphi, \cos \theta)^\top.$$

◇

**2.2.2. Quantum operator postulate.** The quantum operator postulate states that the evolution of a quantum state from  $|\psi\rangle \rightarrow |\psi'\rangle \in \mathcal{H}$  is always achieved via a unitary operator  $U$ , i.e.,

$$(2.11) \quad |\psi'\rangle = U|\psi\rangle, \quad U^\dagger U = I.$$

Here  $U^\dagger$  is the Hermitian conjugate of  $U$ , and  $I$  is the identity map that can be identified with a  $N$ -dimensional identity matrix. The set of all  $N \times N$  unitary matrices is the unitary group, denoted by  $U(N)$ . The set of all  $N \times N$  unitary matrices with determinant 1 forms the special unitary group, denoted by  $SU(N)$ .

This unitary evolution is derived from the system's **Hamiltonian**  $H \in L(\mathcal{H})$ , which is a Hermitian matrix that encapsulates the total energy of the system and thus governs its dynamics. For a time-independent Hamiltonian  $H$ , the state  $|\psi(t)\rangle$  satisfies the Schrödinger equation

$$(2.12) \quad i\partial_t |\psi(t)\rangle = H |\psi(t)\rangle.$$

The corresponding time evolution operator is

$$(2.13) \quad U(t_2, t_1) = e^{-iH(t_2 - t_1)}, \quad \forall t_2 \geq t_1.$$

In particular,  $U(t_2, t_1) = U(t_2 - t_1, 0)$ .



More generally, starting from an initial quantum state  $|\psi(0)\rangle$ , the quantum state can evolve in time, which gives a single parameter family of quantum states denoted by  $\{|\psi(t)\rangle\}$ . These quantum states are related to each other via a quantum evolution operator  $U$ :

$$(2.14) \quad |\psi(t_2)\rangle = U(t_2, t_1) |\psi(t_1)\rangle,$$

where  $U(t_2, t_1)$  is unitary for any given  $t_1, t_2$ . Here  $t_2 > t_1$  refers to quantum evolution forward in time,  $t_2 < t_1$  refers to quantum evolution backward in time, and  $U(t_1, t_1) = I$  for any  $t_1$ .

In quantum computation, a unitary matrix is often referred to as a **quantum gate**.

**Example 2.3.** For a single qubit, the **Pauli matrices** are

$$(2.15) \quad \sigma_x = X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \sigma_y = Y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, \quad \sigma_z = Z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$

Together with the two-dimensional identity matrix, they form a basis of all linear operators on  $\mathbb{C}^2$ .  $\diamond$

Some other commonly used single qubit operators include, to name a few:

- **Hadamard gate**

$$(2.16) \quad H = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

- **Phase gate**

$$(2.17) \quad S = \begin{pmatrix} 1 & 0 \\ 0 & i \end{pmatrix}$$

- **T gate:**

$$(2.18) \quad T = \begin{pmatrix} 1 & 0 \\ 0 & e^{i\pi/4} \end{pmatrix} = \sqrt{S}.$$

When there are notational conflicts, we will use the roman font such as  $H, X$  for these single-qubit gates (for example, to distinguish the Hadamard gate  $H$  from a Hamiltonian  $H$ ). An operator acting on an  $n$ -qubit quantum state space is called an  $n$ -qubit operator.

**Example 2.4.** For  $P \in \{X, Y, Z\}$ , the unitary evolution generated by the Hamiltonian  $H = P$  is a rotation about the corresponding Bloch-sphere axis. Concretely,

$$(2.19) \quad \begin{aligned} R_x(2t) &:= e^{-itX} = \begin{pmatrix} \cos(t) & -i \sin(t) \\ -i \sin(t) & \cos(t) \end{pmatrix}, \\ R_y(2t) &:= e^{-itY} = \begin{pmatrix} \cos(t) & -\sin(t) \\ \sin(t) & \cos(t) \end{pmatrix}, \\ R_z(2t) &:= e^{-itZ} = \begin{pmatrix} e^{-it} & 0 \\ 0 & e^{it} \end{pmatrix}. \end{aligned}$$

For instance, starting from an initial state  $|\psi(0)\rangle = |0\rangle$ , under  $R_x(2t)$  at time  $t = \pi/2$  the state evolves into  $|\psi(\pi/2)\rangle = -i|1\rangle$ , i.e., the  $|1\rangle$  state up to a global phase.  $\diamond$

**THEOREM 2.5** (Spectral theorem of normal matrices). *Given a matrix  $A \in \mathbb{C}^{N \times N}$ , the matrix  $A$  is normal (i.e.,  $A^\dagger A = AA^\dagger$ ), if and only if*

$$(2.20) \quad A = VDV^\dagger.$$

Here,  $D \in \mathbb{C}^{N \times N}$  is a diagonal matrix containing the eigenvalues of  $A$ , and  $V \in U(N)$  is a unitary matrix whose columns are the eigenvectors of  $A$ .

A more general decomposition, which plays a key role throughout the book, is the **singular value decomposition** (SVD).

**THEOREM 2.6** (Singular value decomposition). *Given any matrix  $A \in \mathbb{C}^{M \times N}$ , there exist unitary matrices  $U \in U(M)$  and  $V \in U(N)$ , and a diagonal matrix  $\Sigma \in \mathbb{C}^{M \times N}$  with non-negative real numbers on the diagonal, such that*

$$(2.21) \quad A = U \Sigma V^\dagger.$$

The diagonal entries of  $\Sigma$  are called the *singular values* of  $A$ , the columns of  $U$  are called the *left singular vectors* of  $A$ , and the columns of  $V$  are called the *right singular vectors* of  $A$ .

**Operator exponentials**, also called **matrix exponentials**, gives us a way to express gates as operator exponentials and because the algebra of exponentials makes this representation far easier to work with than explicitly writing the unitary in a matrix representation.

**Definition 2.7** (Matrix function). *For  $A \in \mathbb{C}^{N \times N}$ , and a complex valued function  $f : \mathbb{C} \mapsto \mathbb{C}$ , the matrix function  $f(A)$  is defined as follows:*

- (1) *If  $f$  is an analytic function such that  $f(x) = \sum_{j=0}^{\infty} a_j x^j$  then  $f(A) := \sum_{j=0}^{\infty} a_j A^j$ .*
- (2) *If  $f$  is a complex valued function and  $A$  is a normal matrix such that  $A = V D V^\dagger$  where  $V$  is unitary and  $D := \text{diag}(\lambda_0, \dots, \lambda_{N-1})$  where  $f(\lambda_j) \in \mathbb{C}$ . Then  $f(A) := V f(D) V^\dagger$  where  $f(D) = \text{diag}(f(\lambda_0), \dots, f(\lambda_{N-1}))$ .*

The definition of a matrix exponential can be seen as a direct consequence of either of the above definitions, and both definitions find extensive use in quantum computing. Specifically, using the former definition we have that for any matrix  $A$

$$(2.22) \quad e^A := \sum_{j=0}^{\infty} \frac{A^j}{j!}.$$

Matrix function can also be defined for non-normal matrices using contour integrals (see [Hig08, Chapter 1]).

**Lemma 2.8.** *Let  $A \in \mathbb{C}^{N \times N}$  and let  $U \in U(N)$  be a unitary matrix, then  $U e^A U^\dagger = e^{U A U^\dagger}$ .*

The following result can be viewed as the simplest realization of the **Baker–Campbell–Hausdorff formula** (BCH).

**Lemma 2.9.** *For any  $A, B \in \mathbb{C}^{N \times N}$ , we have*

- (1) *if  $[A, B] = 0$ , then  $e^A e^B = e^{A+B}$ .*
- (2) *if  $[A, [A, B]] = [B, [A, B]] = 0$ , then  $e^A e^B = e^{A+B+\frac{1}{2}[A, B]}$ .*
- (3) *if  $[A, B] \neq 0$ , then  $e^A e^B = e^{A+B+\frac{1}{2}[A, B]} + \mathcal{O}(\max(\|A\|, \|B\|)^3)$ .*

In general, we can express any unitary operator as an exponential of a Hermitian operator. This result is a direct consequence of the definition of the operator exponential.

**Lemma 2.10.** *For any unitary matrix  $U \in U(N)$ , there exists a Hermitian matrix  $H \in \mathbb{C}^{N \times N}$  such that  $U = e^{-iH}$ .*

PROOF. A unitary matrix  $U$  is a normal matrix. According to Theorem 2.5, the unitary matrix  $U$  can be diagonalized as

$$(2.23) \quad U = VDV^\dagger,$$

where  $V \in \mathbb{C}^{N \times N}$  is a unitary matrix and  $D$  is a diagonal matrix. The diagonal entries satisfy  $|D_{ii}| = 1$ . Without loss of generality we can write  $D_{ii} = e^{-i\theta_i}$  where  $\theta_i \in [0, 2\pi)$ . Then define a diagonal matrix  $\Theta_{ii} = \theta_i$ , and  $H = V\Theta V^\dagger$ , we obtain  $U = e^{-iH}$ . Note that the matrix  $H$  is not unique since each  $\theta_i$  can be chosen modulo  $2\pi$ .  $\square$

In many scenarios such as the analysis of quantum simulation using Trotter-Suzuki formulas, we need to find Taylor series expansions of conjugated operators.

**Lemma 2.11.** *Let  $A, B$  be normal matrices in  $\mathbb{C}^{N \times N}$  and let  $t \in \mathbb{R}$ . We then have that*

$$(2.24) \quad e^{At}Be^{-At} = B + \frac{[A, B]t}{1!} + \frac{[A, [A, B]]t^2}{2!} + \frac{[A, [A, [A, B]]]t^3}{3!} + \dots$$

PROOF. We note that the above result is a power series in  $t$ , which must coincide with the Taylor series expansion of the function  $f(t) = e^{At}Be^{-At}$  because the function is analytic. Thus the expression is true if the  $k$ -th derivative of  $f(t)$  at  $t = 0$  is given by the  $k$ -fold commutator. We prove by induction that

$$(2.25) \quad \partial_t^k(e^{At}Be^{-At}) = e^{At}[A, [A, [\dots, [A, B]\dots]]]e^{-At},$$

where the commutator is applied  $k$  times. The base case  $k = 0$  holds trivially. Assume the hypothesis holds for some  $k \geq 0$ . Then

$$(2.26) \quad \begin{aligned} \partial_t^{k+1}(e^{At}Be^{-At}) &= \partial_t(e^{At}[A, [A, [\dots, [A, B]\dots]]]e^{-At}) \\ &= e^{At}(A[A, [\dots, [A, B]\dots]] - [A, [\dots, [A, B]\dots]]A)e^{-At} \\ &= e^{At}[A, [A, [\dots, [A, B]\dots]]]e^{-At}, \end{aligned}$$

where the final expression contains  $k+1$  commutators. This confirms the inductive step. Evaluating at  $t = 0$  yields the coefficients of the Taylor series, completing the proof.  $\square$

**2.2.3. Quantum measurement postulate.** In quantum mechanics, a **quantum observable** is always represented by a Hermitian matrix acting on the state space. The reason for using Hermitian matrices is that they have real eigenvalues, which correspond to the outcome of **quantum measurements**.

A quantum observable  $O \in L(\mathcal{H})$  has the spectral decomposition

$$(2.27) \quad O = \sum_m \lambda_m P_m.$$

Here  $\lambda_m \in \mathbb{R}$  are the eigenvalues of  $O$ , and  $P_m \in L(\mathcal{H})$  is the projection operator onto the eigenspace associated with  $\lambda_m$ . The quantum measurement postulate states that when conducting a measurement on a quantum state  $|\psi\rangle$  with respect to a quantum observable  $O$ , the eigenvalues  $\lambda_m$  represent all the possible results of the measurement. Furthermore, the probability of obtaining a particular outcome  $\lambda_m$  is

$$(2.28) \quad p_m = \langle \psi | P_m | \psi \rangle.$$

Following the measurement, the quantum state collapses to the corresponding eigenspace

$$(2.29) \quad |\psi\rangle \rightarrow \frac{P_m |\psi\rangle}{\sqrt{p_m}}.$$

The set of projection operators satisfies the resolution of identity:

$$(2.30) \quad \sum_m P_m = I.$$

This implies the normalization condition

$$(2.31) \quad \sum_m p_m = \sum_m \langle \psi | P_m | \psi \rangle = \langle \psi | \psi \rangle = 1, \quad \forall |\psi\rangle \in \mathcal{H}.$$

Together with  $p_m \geq 0$ , we find that  $\{p_m\}$  is indeed a probability distribution.

The expectation value of the measurement outcome can be expressed as

$$(2.32) \quad \mathbb{E}_\psi(O) = \sum_m \lambda_m p_m = \sum_m \lambda_m \langle \psi | P_m | \psi \rangle = \left\langle \psi \left| \left( \sum_m \lambda_m P_m \right) \right| \psi \right\rangle = \langle \psi | O | \psi \rangle.$$

**Example 2.12.** Let  $O = X$  be the Pauli  $X$  operator. From the spectral decomposition of  $X$ :

$$(2.33) \quad X |\pm\rangle = \lambda_\pm |\pm\rangle,$$

where  $|\pm\rangle := \frac{1}{\sqrt{2}}(|0\rangle \pm |1\rangle)$ ,  $\lambda_\pm = \pm 1$ , we obtain the eigendecomposition

$$(2.34) \quad O = X = |+\rangle \langle +| - |-\rangle \langle -|.$$

Consider a quantum state  $|\psi\rangle = |0\rangle = \frac{1}{\sqrt{2}}(|+\rangle + |-\rangle)$ , then

$$(2.35) \quad \langle \psi | P_+ | \psi \rangle = \langle \psi | P_- | \psi \rangle = \frac{1}{2}.$$

Therefore the expectation value of the measurement is  $\langle \psi | X | \psi \rangle = 0$ . ◇

**Exercise 2.1.** Prove Eq. (2.34).

#### 2.2.4. Tensor product postulate.

**Definition 2.13** (Tensor product). *The tensor product of two finite dimensional Hilbert spaces  $\mathcal{H}_1$  and  $\mathcal{H}_2$  is a complex vector space, denoted by  $\mathcal{H}_1 \otimes \mathcal{H}_2$ , spanned by vectors of the form  $v \otimes w$  with  $v \in \mathcal{H}_1$  and  $w \in \mathcal{H}_2$ . The bilinear map  $\otimes : \mathcal{H}_1 \times \mathcal{H}_2 \rightarrow \mathcal{H}_1 \otimes \mathcal{H}_2$  satisfies for all  $v, v' \in \mathcal{H}_1$ ,  $w, w' \in \mathcal{H}_2$ , and scalars  $\alpha, \beta \in \mathbb{C}$ :*

- (1)  $(\alpha v + \beta v') \otimes w = \alpha(v \otimes w) + \beta(v' \otimes w)$  and  $v \otimes (\alpha w + \beta w') = \alpha(v \otimes w) + \beta(v \otimes w')$ .
- (2)  $(\alpha v) \otimes w = \alpha(v \otimes w)$  and  $v \otimes (\beta w) = \beta(v \otimes w)$ .

The tensor product is associative in the sense that the two vector spaces  $(\mathcal{H}_1 \otimes \mathcal{H}_2) \otimes \mathcal{H}_3$  and  $\mathcal{H}_1 \otimes (\mathcal{H}_2 \otimes \mathcal{H}_3)$  are isomorphic. Let  $\mathcal{H}_1, \dots, \mathcal{H}_k$  be finite-dimensional Hilbert spaces with inner products  $\langle \cdot | \cdot \rangle_i$  for  $i = 1, 2, \dots, k$ . The tensor product of these  $k$  spaces can be recursively defined as  $\mathcal{H}_1 \otimes \mathcal{H}_2 \otimes \dots \otimes \mathcal{H}_k := \mathcal{H}_1 \otimes (\mathcal{H}_2 \otimes \dots \otimes \mathcal{H}_k)$ , which is spanned by all elements of the form  $v_1 \otimes v_2 \otimes \dots \otimes v_k$  called **product states**, where  $v_i \in \mathcal{H}_i$  for  $i = 1, 2, \dots, k$ . The inner product of two vectors  $v = v_1 \otimes v_2 \otimes \dots \otimes v_k$  and  $w = w_1 \otimes w_2 \otimes \dots \otimes w_k$  in the tensor product space  $\mathcal{H}_1 \otimes \mathcal{H}_2 \otimes \dots \otimes \mathcal{H}_k$  is defined as

$$\langle v | w \rangle = \langle v_1 | w_1 \rangle_1 \cdot \langle v_2 | w_2 \rangle_2 \cdots \langle v_k | w_k \rangle_k.$$

This inner product is extended linearly to the entire tensor product space as

$$\left\langle \sum_i a_i v_i \left| \sum_j b_j w_j \right. \right\rangle = \sum_{i,j} \bar{a}_i b_j \langle v_i | w_j \rangle, \quad a_i, b_j \in \mathbb{C}.$$

The tensor product postulate states that the state space with  $k$  components  $\mathcal{H}_1 \cong \mathbb{C}^{N_1}, \dots, \mathcal{H}_k \cong \mathbb{C}^{N_k}$  is the tensor product of these spaces  $\mathcal{H} = \mathcal{H}_1 \otimes \mathcal{H}_2 \otimes \dots \otimes \mathcal{H}_k$ . Let  $\{|j_i\rangle\}_{j_i \in [N_i]}$  be the basis of  $\mathbb{C}^{N_i}$ , then a general state vector in  $\mathcal{H}$  takes the form

$$(2.36) \quad |\psi\rangle = \sum_{j_1 \in [N_1], \dots, j_k \in [N_k]} \psi_{j_1 \dots j_k} |j_1\rangle \otimes \dots \otimes |j_k\rangle.$$

Here  $\psi_{j_1 \dots j_k} \in \mathbb{C}$  is an entry of a  $k$ -way tensor. Given another state vector

$$(2.37) \quad |\varphi\rangle = \sum_{j_1 \in [N_1], \dots, j_k \in [N_k]} \varphi_{j_1 \dots j_k} |j_1\rangle \otimes \dots \otimes |j_k\rangle,$$

the inner product takes the form

$$(2.38) \quad \langle\psi|\varphi\rangle = \sum_{j_1 \in [N_1], \dots, j_k \in [N_k]} \overline{\psi_{j_1 \dots j_k}} \varphi_{j_1 \dots j_k}.$$

The state space of  $n$ -qubits is  $\mathcal{H} = (\mathbb{C}^2)^{\otimes n} \cong \mathbb{C}^{2^n}$ . We also use a shorthand notation: the tensor product  $\otimes$  may be omitted when the context is clear.

$$(2.39) \quad |01\rangle \equiv |0, 1\rangle \equiv |0\rangle |1\rangle \equiv |0\rangle \otimes |1\rangle, \quad |0^n\rangle \equiv |0^n\rangle \equiv |0\rangle^{\otimes n}.$$

The tensor product operation provides us with a powerful way to describe two independent copies of different vector spaces as a single larger vector space. Further, the tensor product when viewed through this lens does not care about the nature of the form of the Hilbert spaces that are being combined. In fact a particularly important case that we need to consider is the tensor product between two operators.

**Definition 2.14** (Tensor products of linear operators). *Given two finite dimensional Hilbert spaces  $\mathcal{H}_1$  and  $\mathcal{H}_2$ , the tensor product of  $L(\mathcal{H}_1)$  and  $L(\mathcal{H}_2)$ , denoted by  $L(\mathcal{H}_1) \otimes L(\mathcal{H}_2)$ , is a complex vector space spanned by linear operators of the form  $A \otimes B$  with  $A \in L(\mathcal{H}_1)$  and  $B \in L(\mathcal{H}_2)$ . The bilinear map  $\otimes : L(\mathcal{H}_1) \times L(\mathcal{H}_2) \rightarrow L(\mathcal{H}_1) \otimes L(\mathcal{H}_2)$  satisfies for all  $A, B \in L(\mathcal{H}_1)$  and  $C, D \in L(\mathcal{H}_2)$ ,  $v \in \mathcal{H}_1$ ,  $w \in \mathcal{H}_2$  and scalars  $\alpha, \beta \in \mathbb{C}$ :*

- (1)  $(\alpha A + \beta B) \otimes C = \alpha A \otimes C + \beta B \otimes C$  and  $A \otimes (\alpha C + \beta D) = \alpha A \otimes C + \beta A \otimes D$ .
- (2)  $(\alpha A) \otimes B = \alpha A \otimes B = A \otimes (\alpha B)$ .

The space  $L(\mathcal{H}_1) \otimes L(\mathcal{H}_2)$  is isomorphic to  $L(\mathcal{H}_1 \otimes \mathcal{H}_2)$ . The tensor product is also associative in the sense that  $L(\mathcal{H}_1) \otimes (L(\mathcal{H}_2) \otimes L(\mathcal{H}_3))$  is isomorphic to  $(L(\mathcal{H}_1) \otimes L(\mathcal{H}_2)) \otimes L(\mathcal{H}_3)$ . A consequence of this definition is further that the application of multiple tensor products of linear operators on matching tensor products of vectors distributes across the tensor product via

$$(2.40) \quad (A_1 \otimes A_2 \otimes \dots \otimes A_k)(v_1 \otimes v_2 \otimes \dots \otimes v_k) = (A_1 v_1) \otimes (A_2 v_2) \otimes \dots \otimes (A_k v_k).$$

**Example 2.15** (Two qubit system). The state space is  $\mathcal{H} = (\mathbb{C}^2)^{\otimes 2} \cong \mathbb{C}^4$ . The standard basis is (row-major order, i.e., last index is the fastest changing one)

$$(2.41) \quad |00\rangle = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \quad |01\rangle = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \quad |10\rangle = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}, \quad |11\rangle = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}.$$

There are many important quantum operators on the two-qubit quantum system. One of them is the **CNOT gate**, with matrix representation

$$(2.42) \quad \text{CNOT} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

In other words, when acting on the standard basis, we have

$$(2.43) \quad \text{CNOT} \begin{cases} |00\rangle &= |00\rangle \\ |01\rangle &= |01\rangle \\ |10\rangle &= |11\rangle \\ |11\rangle &= |10\rangle \end{cases}.$$

This can be compactly written as

$$(2.44) \quad \text{CNOT} |a\rangle |b\rangle = |a\rangle |a \oplus b\rangle.$$

Here  $a \oplus b = (a + b) \bmod 2$  is the “exclusive or” (XOR) operation.  $\diamond$

**Definition 2.16** (Controlled unitaries). *A controlled unitary operation is a quantum gate that applies a specified unitary operation  $U$  to a set of target qubits only when the control qubits are in a particular state, typically the  $|1\rangle$  state for each control qubit. The single qubit controlled unitary operation can be represented as:*

$$CU = |0\rangle\langle 0| \otimes I + |1\rangle\langle 1| \otimes U.$$

An  $n$ -qubit controlled unitary can be written as:

$$C^n U = (I - |1^n\rangle\langle 1^n|) \otimes I + |1^n\rangle\langle 1^n| \otimes U.$$

The CNOT gate is the same as CX. Controlled unitaries are ubiquitous in quantum algorithms. In particular, it enables conditional logic within quantum circuits.

**Example 2.17** (Multi-qubit Pauli operators). For a  $n$ -qubit quantum system, the Pauli operator acting on the  $i$ -th qubit is denoted by  $P_i$  ( $P = X, Y, Z$ ), i.e.,

$$(2.45) \quad \begin{aligned} X_i &:= I^{\otimes(i-1)} \otimes X \otimes I^{\otimes(n-i)}, \\ Y_i &:= I^{\otimes(i-1)} \otimes Y \otimes I^{\otimes(n-i)}, \\ Z_i &:= I^{\otimes(i-1)} \otimes Z \otimes I^{\otimes(n-i)}. \end{aligned}$$

For example, in a 2-qubit system, following the row-major convention, the matrix representation of  $X_1, X_2$  are

$$(2.46) \quad X_1 = X \otimes I = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}, \quad X_2 = I \otimes X = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

$\diamond$

**Definition 2.18** (Pauli group). *The  $n$ -qubit Pauli group, denoted as  $\mathcal{P}_n$ , is a group that consists of all possible tensor products of  $n$ -qubit Pauli matrices along with multiplicative factors of  $\pm 1$  and  $\pm i$ . Each element of the  $n$ -qubit Pauli group can be represented as*

$$i^k(P_1 \otimes P_2 \otimes \cdots \otimes P_n),$$

where  $k \in \{0, 1, 2, 3\}$ , each  $P_i$  is one of the Pauli matrices  $X, Y, Z$ , or the identity matrix  $I$ , and  $\otimes$  denotes the tensor product.

The  $n$ -qubit Pauli group contains  $4^{n+1}$  elements due to the  $4^n$  possible tensor products of Pauli matrices and identity matrices, each multiplied by one of the four possible phase factors  $\pm 1, \pm i$ . It plays a key role in quantum simulation and quantum error correction. Note that the product of any two elements is another element of the group (up to a phase factor), and every element is its own inverse (up to a phase factor).

**Definition 2.19** (Clifford group). *The  $n$ -qubit Clifford group, denoted as  $\mathcal{C}_n$ , is a group of unitary operators that normalizes the  $n$ -qubit Pauli group  $\mathcal{P}_n$ . This means that for every Clifford operator  $C \in \mathcal{C}_n$  and every Pauli operator  $P \in \mathcal{P}_n$ , there exists a Pauli operator  $P' \in \mathcal{P}_n$  such that*

$$CPC^\dagger = P'.$$

The Clifford group includes all elements of the Pauli group, the Hadamard gate  $H$ , the phase gate  $S$ , and the CNOT gate. It can be generated by  $\{H, S, \text{CNOT}\}$ .

**Example 2.20.** The single-qubit Pauli group  $\mathcal{P}_1$  is defined as the group generated by the Pauli matrices  $X, Y, Z$  together with the phase factor  $i$ :

$$\mathcal{P}_1 = \{i^k P \mid k \in \{0, 1, 2, 3\}, P \in \{I, X, Y, Z\}\}.$$

We show that  $\mathcal{P}_1$  can be generated by the set  $\{H, S\}$ . First, we obtain the Pauli  $Z$  operator by squaring the phase gate:

$$S^2 = \begin{pmatrix} 1 & 0 \\ 0 & i \end{pmatrix}^2 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} = Z.$$

Next, we utilize the property that the Hadamard gate transforms  $Z$  into  $X$  under conjugation. Since  $H$  is Hermitian and unitary, we have:

$$X = HZH = HS^2H.$$

The Pauli  $Y$  operator can be generated by conjugating  $X$  by  $S$ . We compute the conjugate transpose  $S^\dagger = \text{diag}(1, -i)$  and verify the relation:

$$\begin{aligned} SXS^\dagger &= \begin{pmatrix} 1 & 0 \\ 0 & i \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & -i \end{pmatrix} \\ &= \begin{pmatrix} 0 & 1 \\ i & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & -i \end{pmatrix} = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} = Y. \end{aligned}$$

Since  $S$  is unitary and  $S^4 = I$ , we have  $S^\dagger = S^{-1} = S^3$ . Thus,  $Y = SXS^3$ .

Finally, since  $XYZ = iI$ , we conclude that  $\{H, S\}$  generates the entire Pauli group  $\mathcal{P}_1$ .

Since

$$T^2 = \begin{pmatrix} 1 & 0 \\ 0 & e^{i\pi/4} \end{pmatrix}^2 = \begin{pmatrix} 1 & 0 \\ 0 & e^{i\pi/2} \end{pmatrix} = S,$$

it immediately follows that  $\{H, T\}$  also generates  $\mathcal{P}_1$ . ◇

The Clifford group plays an important role in many areas. In quantum error correction, Clifford operations can transform certain errors into forms that are more easily correctable. This makes it desirable to choose Clifford gates to be part of a universal gate set (the most common one is Clifford +  $T$ ). Additionally, the Gottesman–Knill theorem states that any quantum circuit using only Clifford gates on computational basis states and measurements in the computational basis can be efficiently simulated classically.

**Example 2.21.** We can concisely describe block matrices stored within a larger matrix. The matrix representation of  $T \in L(\mathbb{C}^{NM} \otimes \mathbb{C}^{NM})$ , when writing in the block form,

$$(2.47) \quad T = \begin{pmatrix} T_{0,0} & \cdots & T_{0,N-1} \\ \vdots & \ddots & \vdots \\ T_{N-1,0} & \cdots & T_{N-1,N-1} \end{pmatrix}, \quad T_{ij} \in \mathbb{C}^{M \times M}.$$

can be rewritten as

$$(2.48) \quad T = \sum_{i,j \in [N]} |e_i\rangle\langle e_j| \otimes T_{ij}.$$

◇

The notation for partial inner products and partial applications of operators is used throughout this book, particularly in the context of block-encoding.

**Definition 2.22** (Partial inner product). *Consider two finite dimensional Hilbert spaces  $\mathcal{H}_A \cong \mathbb{C}^N$  with an orthonormal basis  $\{|e_i\rangle\}_{i \in [N]}$ , and  $\mathcal{H}_B \cong \mathbb{C}^M$  with an orthonormal basis  $\{|f_i\rangle\}_{i \in [M]}$ . The partial inner product  $\langle \cdot | \cdot \rangle$  is a map  $\mathcal{H}_A \times (\mathcal{H}_A \otimes \mathcal{H}_B) \rightarrow \mathcal{H}_B$  defined as follows. For any  $v \in \mathcal{H}_A$ ,  $w \in \mathcal{H}_A \otimes \mathcal{H}_B$*

$$(2.49) \quad \langle v | w \rangle = \sum_{ij} (\langle v | e_i \rangle \langle e_i, f_j | w \rangle) |f_j\rangle \in \mathcal{H}_B.$$

*With some abuse of notation, the partial inner product  $\langle \cdot | \cdot \rangle$  also denotes a map:  $(\mathcal{H}_A \otimes \mathcal{H}_B) \times \mathcal{H}_A \rightarrow \mathcal{H}_B^*$  according to*

$$(2.50) \quad \langle w | v \rangle = \sum_{ij} (\langle e_i | v \rangle \langle w | e_i, f_j \rangle) \langle f_j | \in \mathcal{H}_B^*.$$

This definition of a partial inner product has been used in the literature in several works such as [LC17]. A problem with the notation though is that it requires that the reader pay close attention to the dimensions of the objects in question in order to infer the dimension of the output with a partial inner product. This runs counter to the advantages of Dirac notation which can be confusing when used in the context of conventional Dirac notation where the inner product is always a scalar. While its brevity is an advantage, great care must be taken when using the above notation to avoid making mistakes about the shape of the output.

**Example 2.23.** Let  $|v\rangle = \frac{1}{\sqrt{2}}|0\rangle + \frac{1}{\sqrt{2}}|1\rangle$  be a one-qubit state,  $|w\rangle = |0\rangle \otimes (|00\rangle + |11\rangle) + |1\rangle \otimes (|01\rangle + |10\rangle)$  be a three-qubit state, then the partial inner product

$$(2.51) \quad \langle v | w \rangle = \frac{1}{\sqrt{2}}(|00\rangle + |11\rangle) + \frac{1}{\sqrt{2}}(|01\rangle + |10\rangle)$$

is a two-qubit state.

◇



**Example 2.24.** Let  $w = \sum_{i \in [N]} |e_i\rangle \otimes |w_i\rangle$  be reshaped into a matrix

$$(2.52) \quad W = \begin{pmatrix} w_0 & \cdots & w_{N-1} \end{pmatrix} \in \mathbb{C}^{N \times M}.$$

Then the partial inner product  $\langle e_i | w \rangle$  for  $i \in [N]$  picks out the  $i$ -th column  $w_i$ . Similarly, the partial inner product  $\langle w | e_i \rangle$  picks out the  $i$ -th row of  $W^\dagger$ , which is  $w_i^\dagger$ .  $\diamond$

The partial inner product between pure states provides a natural way to focus our attention on one of the subspaces involved. Sometimes however, we will wish to apply a transformation on the system in question. This generalizes the concept of the partial inner product, and will be vital in our later discussion on block-encoding in Chapter 9.

**Definition 2.25** (Partial application of operators). *Consider two finite dimensional Hilbert spaces  $\mathcal{H}_A \cong \mathbb{C}^N$  with an orthonormal basis  $\{|e_i\rangle\}_{i \in [N]}$ , and  $\mathcal{H}_B \cong \mathbb{C}^M$  with an orthonormal basis  $\{|f_i\rangle\}_{i \in [M]}$ . A partial application is a map  $(\mathcal{H}_A^* \otimes L(\mathcal{H}_B)) \times (\mathcal{H}_A \otimes \mathcal{H}_B) \rightarrow \mathcal{H}_B$  so that for  $|v\rangle \in \mathcal{H}_A$ ,  $C \in L(\mathcal{H}_B)$ ,  $|u\rangle \in \mathcal{H}_A \otimes \mathcal{H}_B$ ,*

$$(2.53) \quad (\langle v | \otimes C) |u\rangle = \sum_{jk} (\langle v | e_j \rangle \langle e_j, f_k | u \rangle) (C | f_k \rangle) \in \mathcal{H}_B.$$

Similarly we define

$$(2.54) \quad \langle u | (|v\rangle \otimes C) = \sum_{jk} (\langle e_j | v \rangle \langle u | e_j, f_k \rangle) (\langle f_k | C) \in \mathcal{H}_B^*.$$

**Example 2.26.** The partial inner product can also be viewed as a partial application of the identity gate, i.e., for  $|v\rangle \in \mathcal{H}_A$ ,  $I \in L(\mathcal{H}_B)$ ,  $|u\rangle \in \mathcal{H}_A \otimes \mathcal{H}_B$

$$(2.55) \quad \langle v | u \rangle = (\langle v | \otimes I) |u\rangle, \quad \langle u | v \rangle = \langle u | (|v\rangle \otimes I).$$

For  $T = \sum_{jk} |e_j\rangle\langle e_k| \otimes T_{jk}$ , the quantity  $(\langle e_i | \otimes I)T$  can be represented as a rectangular matrix that consists of the  $i$ -th block row of  $T$ :

$$(2.56) \quad \begin{aligned} (\langle e_i | \otimes I)T &= (\langle e_i | \otimes I) \sum_{jk} |e_j\rangle\langle e_k| \otimes T_{jk} \\ &= \sum_k \langle e_k | \otimes T_{ik} \equiv (T_{i,0} \quad \cdots \quad T_{i,N-1}), \end{aligned}$$

Similarly,  $T(|e_j\rangle \otimes I)$  picks out the  $j$ -th block column of the matrix  $T$ .

$$(2.57) \quad \begin{aligned} T(|e_j\rangle \otimes I) &= \sum_{ik} \delta_{jk} |e_i\rangle \otimes T_{ik} \\ &= \sum_i |e_i\rangle \otimes T_{ij} \equiv \begin{pmatrix} T_{0,j} \\ \vdots \\ T_{N-1,j} \end{pmatrix}, \end{aligned}$$

and  $(\langle e_i | \otimes I)T(|e_j\rangle \otimes I)$  returns the  $(i, j)$ -th block  $T_{ij}$  as can be seen via

$$(2.58) \quad (\langle e_i | \otimes I)T(|e_j\rangle \otimes I) = (\langle e_i | \otimes I) \sum_k |e_k\rangle \otimes T_{kj} = T_{ij}.$$

With some abuse of notation, we may omit the  $\otimes I$  notation, so  $(\langle e_i | \otimes I)T(|e_j\rangle \otimes I)$  may be written simply as  $\langle e_i | T | e_j \rangle$ .  $\diamond$

### 2.3. Density operator

So far all quantum states encountered have been described by a single state vector  $|\psi\rangle$ . How to describe a classical mixture of state vectors, such as the state after a measurement process? How can the state of a subsystem within a larger quantum system be defined? The answer to these questions requires the formulation of the density operator (also called density matrix).

**Definition 2.27** (Density operator). *A linear operator  $\rho \in L(\mathbb{C}^N)$  is called a density operator, if  $\rho \succeq 0$ , and  $\text{Tr } \rho = 1$ . The set of all density operators is denoted by  $\mathcal{D}(\mathbb{C}^N)$ .*

The density operator corresponding to a state vector  $|\psi\rangle$  is a rank-1 matrix

$$(2.59) \quad \rho = |\psi\rangle\langle\psi|.$$

Recall that quantum mechanics postulates that  $|\psi\rangle$  and  $|\psi'\rangle = e^{i\theta} |\psi\rangle$  represent the same physical state. This statement is more natural from the perspective of the density operator, since

$$(2.60) \quad \rho' = |\psi'\rangle\langle\psi'| = e^{i\theta} |\psi\rangle e^{-i\theta} \langle\psi| = \rho.$$

In physics, such an irrelevant phase factor is referred to as a gauge degree of freedom. The density operator  $\rho$  encapsulates the same physical information as is present in  $|\psi\rangle$ , but with the added benefit of being invariant to the gauge choice.

With some abuse of terminology, throughout this book, both the density operator  $\rho$  and the state vector  $|\psi\rangle$  are called quantum states. A rank-1 density operator is called a **pure state**.

**Exercise 2.2.** Prove that all eigenvalues of a density operator  $\rho$  belong to  $[0, 1]$ . Furthermore,  $\rho^2 \preceq \rho$ , and the equality holds if and only if  $\rho$  is a pure state.

If  $\rho$  is not a pure state, then it is called a **mixed state**. We can diagonalize the density matrix as

$$(2.61) \quad \rho = \sum_i p_i |\psi_i\rangle\langle\psi_i| =: \sum_i p_i \rho_i,$$

where all state vectors  $|\psi_i\rangle$  are orthogonal to each other, and each  $\rho_i$  is a pure state. On the other hand, if we have the ability to prepare each pure state  $\rho_i$ , then to create the mixed state  $\rho$ , all we need to do is prepare a state  $\rho_i$  randomly, with the probability of preparing each state given by  $p_i$ . In essence, a mixed state can be seen as a **classical ensemble** of pure quantum states. In particular, an  $n$ -qubit state  $\rho = \frac{I}{2^n}$  is called the **maximally mixed state**.

Let  $\{\rho_j\}$  be a set of density operators. With any discrete probability distribution  $\{p_j\}$ , define  $\rho' = \sum_j p_j \rho_j$ . Then  $\rho' \succeq 0$  and  $\text{Tr}[\rho'] = \sum_j p_j \text{Tr}[\rho_j] = \sum_j p_j = 1$ . Therefore  $\rho'$  is a density operator. In other words, a classical ensemble of (pure or mixed) density operators is also a density operator.

**Example 2.28** (Expectation value of a quantum observable). Let us consider the expectation value of an observable  $O$  with respect to a mixed state  $\rho$ . Since the expectation value with respect to a pure state is

$$(2.62) \quad \langle O \rangle_{\rho_i} = \langle \psi_i | O | \psi_i \rangle = \text{Tr}[O \rho_i],$$

if we obtain the expectation value for a mixed state that obtains a pure state  $\rho_i$  with probability  $p_i$ , the expectation value is concisely written as

$$(2.63) \quad \langle O \rangle_{\rho} = \sum_i p_i \text{Tr}[O \rho_i] = \text{Tr}[O \rho].$$

◇

The measurement process can be described without referring to a quantum observable. A quantum measurement can be described by a set of measurement operators  $\{M_m\}$ , where  $m$  labels the different possible outcomes of the measurement. The operators  $M_m$  act on the state space  $\mathcal{H}$  of the system and satisfy the completeness relation:  $\sum_m M_m^\dagger M_m = I$ . After a measurement described by  $M_m$  is made on a quantum system in a state  $\rho$ , the probability of getting result  $m$  is given by

$$(2.64) \quad p_m = \text{Tr}[M_m \rho M_m^\dagger].$$

If outcome  $m$  occurs, then the state of the quantum system collapses to a new state

$$(2.65) \quad \rho'_m = \frac{M_m \rho M_m^\dagger}{\text{Tr}[M_m \rho M_m^\dagger]}.$$

The density operator of the resulting ensemble is

$$(2.66) \quad \rho' = \sum_m p_m \rho'_m = \sum_m M_m \rho M_m^\dagger.$$

If each  $M_m$  is a projection operator denoted by  $P_m$ , then  $\{P_m\}$  is called a **projective measurement**. When a quantum observable is measured, the action that is performed on the quantum system is a projective measurement. That is, the state of the system is projected onto an eigenstate of the observable, corresponding to the obtained result of the measurement.

**Example 2.29** (Projective measurement). Let the initial state  $\rho = |\psi\rangle\langle\psi|$  be a pure state subject to a projective measurement  $\{P_m\}_m$ . After measurement, the system collapses into a state  $|\psi_m\rangle = P_m |\psi\rangle / \sqrt{p_m}$  with probability  $p_m = \langle\psi|P_m|\psi\rangle$ . If we attempt to represent it by a pure state, one natural choice seems to be  $|\psi'\rangle = \sum_m \sqrt{p_m} |\psi_m\rangle$ . However, using the normalization condition of the projective measurement in Eq. (2.30)

$$(2.67) \quad \sum_m \sqrt{p_m} |\psi_m\rangle = \sum_m \sqrt{p_m} P_m |\psi\rangle / \sqrt{p_m} = \sum_m P_m |\psi\rangle = |\psi\rangle.$$

In other words, state before and after the measurement is exactly the same! This clearly does not make sense.

Instead, the resulting state should be represented by a mixed state

$$(2.68) \quad \rho' = \sum_m p_m |\psi_m\rangle\langle\psi_m| = \sum_m P_m |\psi\rangle\langle\psi| P_m = \sum_m P_m \rho P_m.$$

◇

The partial trace is an operation on a joint quantum state (often representing a composite system), which effectively “traces out” one or more subsystems to leave a reduced density operator for the remaining subsystem(s). The operation is widely used in quantum mechanics, especially in the study of open quantum systems, quantum information, and quantum computation.

**Definition 2.30** (Partial trace). Consider two finite dimensional Hilbert spaces  $\mathcal{H}_A \cong \mathbb{C}^N$  with an orthonormal basis  $\{|e_i\rangle\}_{i \in [N]}$ , and  $\mathcal{H}_B \cong \mathbb{C}^M$ , and  $T \in L(\mathcal{H}_A \otimes \mathcal{H}_B)$ . The partial trace over  $\mathcal{H}_A$ , denoted by  $\text{Tr}_A(T)$  is an element in  $L(\mathcal{H}_B)$  defined as:

$$(2.69) \quad \text{Tr}_A(T) = \sum_{i \in [N]} (\langle e_i | \otimes I) T (|e_i\rangle \otimes I).$$

The partial trace  $\text{Tr}_B(T)$  is defined similarly.

**Example 2.31.** The matrix representation of  $T \in L(\mathbb{C}^N \otimes \mathbb{C}^M)$  takes the form of a block matrix

$$(2.70) \quad T = \begin{pmatrix} T_{0,0} & \cdots & T_{0,N-1} \\ \vdots & \ddots & \vdots \\ T_{N-1,0} & \cdots & T_{N-1,N-1} \end{pmatrix}, \quad T_{ij} \in \mathbb{C}^{M \times M}.$$

Then

$$(2.71) \quad \text{Tr}_A(T) = \sum_{i \in [N]} T_{ii}$$

is the sum of all diagonal blocks. ◇

Given a density operator  $\rho \in \mathcal{D}(\mathcal{H}_A \otimes \mathcal{H}_B)$ , the partial trace

$$(2.72) \quad \rho_A = \text{Tr}_B[\rho] \in \mathcal{D}(\mathcal{H}_A), \quad \rho_B = \text{Tr}_A[\rho] \in \mathcal{D}(\mathcal{H}_B)$$

are called **reduced density operators**. In particular, if  $\rho = \rho_1 \otimes \rho_2$ , then

$$(2.73) \quad \text{Tr}_B[\rho] = \rho_1, \quad \text{Tr}_A[\rho] = \rho_2.$$

Note that even if  $\rho$  is a pure state, in general, the reduced density operators  $\rho_A, \rho_B$  are mixed states.

If a quantum observable is defined only on the subsystem  $A$ , i.e.,  $O = O_A \otimes I_B$  and  $O_A = \sum_m \lambda_m P_m$ , then when measuring a quantum state  $\rho$  with respect to  $O$ , the probability of obtaining  $\lambda_m$ , and the expectation value only depend on the reduced density matrix  $\rho_A$ :

$$(2.74) \quad p_m = \text{Tr}[(P_m \otimes I)\rho] = \text{Tr}[P_m \text{Tr}_B[\rho]] = \text{Tr}[P_m \rho_A], \quad \mathbb{E}_\rho[O] = \text{Tr}[(O_A \otimes I)\rho] = \text{Tr}[O_A \rho_A].$$

**Exercise 2.3.** The **Bell state** (also called the EPR pair) is defined to be

$$(2.75) \quad |\psi\rangle = \frac{1}{\sqrt{2}}(|00\rangle + |11\rangle) = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \end{pmatrix}.$$

Use the partial trace over the second qubit to prove that the Bell state cannot be written as any product state  $|a\rangle \otimes |b\rangle$ .

**Example 2.32** (Purification of mixed state). Any mixed state can always be dilated to a pure state using ancilla qubits. In particular, any  $n$ -qubit mixed state  $\rho$  can be expressed as  $\sum_j p_j |\lambda_j\rangle\langle\lambda_j|$  where  $|\lambda_j\rangle$  are the eigenvectors of  $\rho$ , and  $p_j$  is the corresponding eigenvalue. Given this we can construct a  $2n$ -qubit pure state

$$(2.76) \quad |\rho\rangle := \sum_j \sqrt{p_j} |\lambda_j\rangle_A |\lambda_j\rangle_B.$$

Then  $\text{Tr}_B(|\rho\rangle\langle\rho|) = \rho$ . ◇

A more general concept than projective measurement is called **generalized measurement**, also called positive operator-valued measure (POVM).

**Definition 2.33.** A **positive operator-valued measure** (POVM) is a set of positive semidefinite operators  $\{E_m\}$  that sum to the identity:

$$(2.77) \quad \sum_m E_m = I, \quad E_m \succeq 0.$$

If a quantum system is in state  $\rho$ , the probability of obtaining outcome  $m$  is given by

$$(2.78) \quad p_m = \text{Tr}[E_m \rho].$$

Unlike projective measurements, the elements  $E_m$  of a POVM are not necessarily orthogonal, nor are they required to be projection operators (i.e.,  $E_m^2$  need not equal  $E_m$ ). However, POVMs provide the most general description of quantum measurements. On the other hand, the Naimark's dilation theorem (see e.g. [Wat18, Chapter 2.3]) tells us that any generalized measurement can be implemented by coupling the system of interest to an ancilla system and performing a standard projective measurement on the composite system.

**THEOREM 2.34** (Naimark's dilation theorem). *Every POVM can be realized as a projective measurement on a larger Hilbert space. Specifically, given a POVM  $\{E_m\}$  on  $\mathcal{H}_A$ , there exists an auxiliary Hilbert space  $\mathcal{H}_B$ , a pure state  $|0\rangle_B \in \mathcal{H}_B$ , and a projective measurement  $\{P_m\}$  on  $\mathcal{H}_A \otimes \mathcal{H}_B$  such that for any state  $\rho$  on  $\mathcal{H}_A$ :*

$$(2.79) \quad \text{Tr}[E_m \rho] = \text{Tr}[P_m(\rho \otimes |0\rangle\langle 0|_B)].$$

**PROOF.** Since each  $E_m$  is positive semidefinite, we can define  $M_m = \sqrt{E_m}$  such that  $M_m^\dagger M_m = E_m$ . Let  $\mathcal{H}_B$  be a Hilbert space with an orthonormal basis  $\{|m\rangle\}$  corresponding to the indices of the POVM elements. We define a linear operator  $V : \mathcal{H}_A \rightarrow \mathcal{H}_A \otimes \mathcal{H}_B$  by its action on an arbitrary state  $|\psi\rangle \in \mathcal{H}_A$ :

$$(2.80) \quad V|\psi\rangle = \sum_m M_m |\psi\rangle \otimes |m\rangle_B.$$

This operator is an isometry because

$$(2.81) \quad \langle V\psi | V\psi \rangle = \sum_{m,n} \langle \psi | M_m^\dagger M_n | \psi \rangle \langle m | n \rangle = \sum_m \langle \psi | M_m^\dagger M_m | \psi \rangle = \langle \psi | \left( \sum_m E_m \right) | \psi \rangle = \langle \psi | \psi \rangle.$$

We can extend this isometry to a unitary operator  $U$  acting on  $\mathcal{H}_A \otimes \mathcal{H}_B$  such that  $U(|\psi\rangle \otimes |0\rangle_B) = V|\psi\rangle$ . Now, define the projective measurement on the composite system by the projectors  $\Pi_m = I_A \otimes |m\rangle\langle m|_B$ . Let  $P_m = U^\dagger \Pi_m U$ . Since  $U$  is unitary and  $\{\Pi_m\}$  are orthogonal projectors summing to identity,  $\{P_m\}$  is a valid projective measurement. Finally, we verify the probability condition:

$$(2.82) \quad \begin{aligned} \text{Tr}[P_m(\rho \otimes |0\rangle\langle 0|_B)] &= \text{Tr}[U^\dagger \Pi_m U(\rho \otimes |0\rangle\langle 0|_B)] \\ &= \text{Tr}[\Pi_m U(\rho \otimes |0\rangle\langle 0|_B)U^\dagger] \\ &= \text{Tr}[(I_A \otimes |m\rangle\langle m|_B)V\rho V^\dagger]. \end{aligned}$$

Using the definition of  $V$ , we have  $V\rho V^\dagger = \sum_{k,l} M_k \rho M_l^\dagger \otimes |k\rangle\langle l|_B$ . Substituting this back,

$$(2.83) \quad \begin{aligned} \text{Tr}[P_m(\rho \otimes |0\rangle\langle 0|_B)] &= \text{Tr} \left[ (I_A \otimes |m\rangle\langle m|_B) \sum_{k,l} M_k \rho M_l^\dagger \otimes |k\rangle\langle l|_B \right] \\ &= \text{Tr}[M_m \rho M_m^\dagger] = \text{Tr}[M_m^\dagger M_m \rho] = \text{Tr}[E_m \rho]. \end{aligned}$$

□

### 2.4. Quantum circuit

Nearly all quantum algorithms operate on multi-qubit quantum systems. When quantum operators operate on two or more qubits, writing down quantum states in terms of its components as in Eq. (2.36) quickly becomes cumbersome. The language of **quantum circuit** offers a graphical and compact manner for writing down the procedure of applying a sequence of quantum operators to a quantum state. For more details see [NC00, Section 4.2, 4.3].

In the quantum circuit language, time flows from the left to right, i.e., the input quantum state appears on the left, and the quantum operator appears on the right, and each “wire” represents a qubit i.e.,

$$|\psi\rangle \text{ --- } \boxed{U} \text{ --- } U|\psi\rangle$$

Here are a few examples:

$$|0\rangle \text{ --- } \boxed{X} \text{ --- } |1\rangle \quad |1\rangle \text{ --- } \boxed{Z} \text{ --- } -|1\rangle \quad |0\rangle \text{ --- } \boxed{H} \text{ --- } |+\rangle$$

which is a graphical way of writing

$$(2.84) \quad X|0\rangle = |1\rangle, \quad Z|1\rangle = -|1\rangle, \quad H|0\rangle = |+\rangle.$$

The relation between these states can be expressed in terms of the following diagram

$$(2.85) \quad \begin{array}{ccc} |0\rangle & \xrightarrow{X} & |1\rangle \\ \downarrow H & & \downarrow H \\ |+\rangle & \xrightarrow{Z} & |-\rangle \end{array}$$

Also verify that

$$\begin{array}{ccc} |0\rangle & \text{--- } \boxed{X} \text{ ---} & |1\rangle \\ |0\rangle & \text{---} & |0\rangle \end{array}$$

which is a graphical way of writing

$$(2.86) \quad (X \otimes I)|00\rangle = |10\rangle.$$

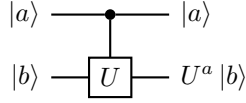
Note that the input state can be general, and in particular does not need to be a product state. For example, if the input is a Bell state (2.75), we just apply the quantum operator to  $|00\rangle$  and  $|11\rangle$ , respectively and multiply the results by  $1/\sqrt{2}$  and add together. To distinguish with other symbols, these single qubit gates may be either written as  $X, Y, Z, H$  or (using the roman font)  $X, Y, Z, H$ .

The quantum circuit for the CNOT gate is

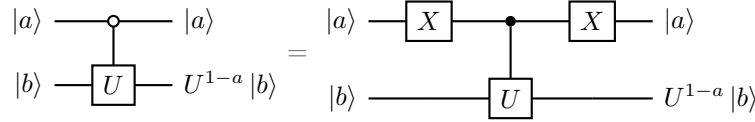
$$\begin{array}{ccc} |a\rangle & \text{---} \bullet \text{---} & |a\rangle \\ & \downarrow \oplus & \\ |b\rangle & \text{---} \oplus \text{---} & |a \oplus b\rangle \end{array}$$

Here the “dot” means that the quantum gate connected to the dot only becomes active if the state of the qubit 0 (called the control qubit) is  $a = 1$ . This justifies the name of the CNOT gate (controlled NOT).

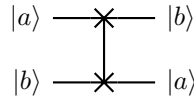
Similarly,



is the controlled  $U$  gate for some unitary  $U$ . Here  $U^a = I$  if  $a = 0$ . The CNOT gate can be obtained by setting  $U = X$ . Sometimes we want to control a unitary only if the control qubit is zero rather than 1. In this case, we represent the control using a hollow circle as shown below.

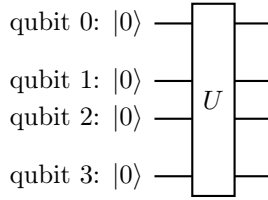


Another commonly used two-qubit gate is the **SWAP gate**, which swaps the state in the 0-th and the 1-st qubits.



**Exercise 2.4.** Write down the matrix representation of the SWAP gate.

Quantum operators applied to multiple qubits can be written in a similar manner:

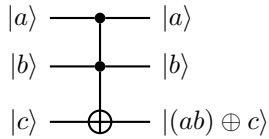


For a multi-qubit quantum circuit, unless stated otherwise, the first qubit will be referred to as the qubit 0, and the second qubit as the qubit 1, etc.

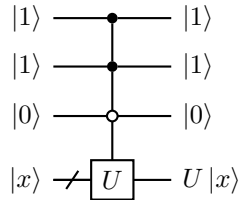
When the context is clear, we may also use a more compact notation for the multi-qubit quantum operators:

$$|0\rangle^{\otimes 4} \text{ --- } \boxed{U} \text{ ---} \Leftrightarrow |0\rangle^{\otimes 4} \equiv \boxed{U} \equiv \Leftrightarrow |0\rangle^{\otimes 4} \text{ --- } \boxed{U} \text{ ---}$$

One useful multiple qubit gate is the **Toffoli gate** (or controlled-controlled-NOT, CCNOT gate).



We may also want to apply a  $n$ -qubit unitary  $U$  only when certain conditions are met

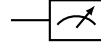


where the empty circle means that the gate being controlled only becomes active when the value of the control qubit is 0. This can be used to write down the quantum “if” statements, i.e., when the qubits 0, 1 are at the  $|1\rangle$  state and the qubit 2 is at the  $|0\rangle$  state, then apply  $U$  to  $|x\rangle$ .

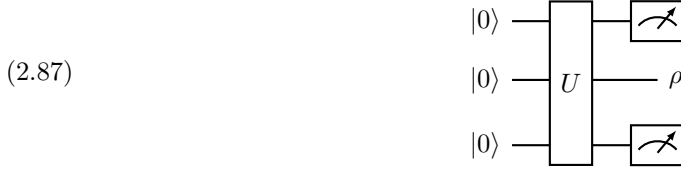
A set of qubits is often called a **quantum register** (or register for short). For example, in the picture above, the main quantum state of interest (an  $n$  qubit quantum state  $|x\rangle$ ) is called the system register. The first 3 qubits can be called the control register. When multiple registers are present, we can distinguish them by writing  $|x\rangle_A |y\rangle_B$ , so that we can refer to the quantum state associated with the qubits in registers  $A$  and  $B$ , respectively.

In quantum computation, a classical bit-string is denoted as  $x \in \{0, 1\}^n$ , and the corresponding  $|x\rangle$  is called a **classical state**. The set of all classical states form the **computational basis** of an  $n$ -qubit system. It is worth noting that  $\{|x\rangle \langle x| \mid x \in \{0, 1\}^n\}$  forms a set of projective measurement operators, which can be identified with the simultaneous measurement with respect to Pauli-Z operators  $Z_1, \dots, Z_n$ . Consequently, when a measurement is performed with respect to the Pauli-Z operator, it is called a measurement in the computational basis.

The circuit symbol for the quantum measurement with respect to a single Pauli-Z is



**Example 2.35** (Measure Pauli-Z operators). For a quantum state  $|\psi\rangle$ , the measurement of a multi-qubit Pauli-Z operator of the form  $(Z_1)^{a_1} \dots (Z_n)^{a_n}$ , where  $a_1, \dots, a_n \in \{0, 1\}$  can be directly implemented at the circuit level. For example, for a 3-qubit system, the following circuit



measures the outcome of  $Z_1$  and  $Z_3$ , yielding 4 possible outcomes  $\{00, 01, 10, 11\}$  with respective probabilities  $\{p(00), p(01), p(10), p(11)\}$ . Now consider an observable  $O = Z_1 Z_3$  whose eigenvalues are 1 and  $-1$ . The probability of obtaining each eigenvalue is

(2.88) 
$$p(O = 1) = p(00) + p(11), \quad p(O = -1) = p(01) + p(10).$$

◇

**Example 2.36** (Hadamard test circuit). The Hadamard test is a useful tool for computing the expectation value of an unitary operator with respect to a state, i.e.,  $\langle \psi | U | \psi \rangle$ . It can be used to solve the phase estimation problem. The Hadamard test uses two circuits to estimate the real and imaginary part of the expectation value separately.

The (real) Hadamard test is the quantum circuit in Fig. 2.1 for estimating  $\text{Re} \langle \psi | U | \psi \rangle$ .

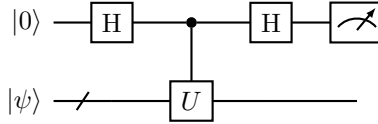


FIGURE 2.1. Hadamard test for  $\text{Re} \langle \psi | U | \psi \rangle$ .



To verify this, we find that the circuit transforms  $|0\rangle |\psi\rangle$  as

$$\begin{aligned} |0\rangle |\psi\rangle &\xrightarrow{H \otimes I} \frac{1}{\sqrt{2}}(|0\rangle + |1\rangle) |\psi\rangle \\ &\xrightarrow{c-U} \frac{1}{\sqrt{2}}(|0\rangle |\psi\rangle + |1\rangle U |\psi\rangle) \\ &\xrightarrow{H \otimes I} \frac{1}{2} |0\rangle (|\psi\rangle + U |\psi\rangle) + \frac{1}{2} |1\rangle (|\psi\rangle - U |\psi\rangle). \end{aligned}$$

The probability of measuring the qubit 0 to be in state  $|0\rangle$  is

$$(2.89) \quad p(0) = \frac{1}{2}(1 + \text{Re} \langle \psi | U | \psi \rangle).$$

This is well defined since  $-1 \leq \text{Re} \langle \psi | U | \psi \rangle \leq 1$ .

To obtain the imaginary part, we can use the circuit in Fig. 2.2 called the (imaginary) Hadamard test, where

$$(2.90) \quad S = \begin{pmatrix} 1 & 0 \\ 0 & i \end{pmatrix}$$

is called the phase gate.

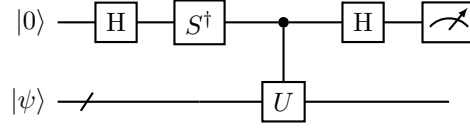


FIGURE 2.2. Hadamard test for  $\text{Im} \langle \psi | U | \psi \rangle$ .

Similar calculation shows the circuit transforms  $|0\rangle |\psi\rangle$  to the state

$$(2.91) \quad \frac{1}{2} |0\rangle (|\psi\rangle - iU |\psi\rangle) + \frac{1}{2} |1\rangle (|\psi\rangle + iU |\psi\rangle).$$

Therefore the probability of measuring the qubit 0 to be in state  $|0\rangle$  is

$$(2.92) \quad p(0) = \frac{1}{2}(1 + \text{Im} \langle \psi | U | \psi \rangle).$$

Combining the results from the two circuits, we obtain the estimate to  $\langle \psi | U | \psi \rangle$ .

◇

**Example 2.37** (Overlap estimate using the SWAP test). A special case of the Hadamard test is called the **SWAP test**, which can be used to estimate the overlap of two quantum states  $|\langle \varphi | \psi \rangle|$ . The quantum circuit for the swap test is

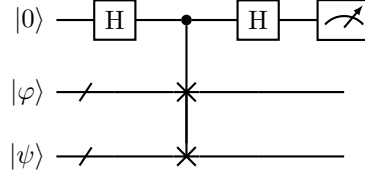


FIGURE 2.3. Circuit for the SWAP test.

Note that this is exactly the Hadamard test with  $U$  being the swap gate. Direct calculation shows that the probability of measuring the first qubit and obtaining outcome 0 is

$$(2.93) \quad p(0) = \frac{1}{2}(1 + \operatorname{Re} \langle \varphi, \psi | \psi, \varphi \rangle) = \frac{1}{2}(1 + |\langle \varphi | \psi \rangle|^2).$$

◇

### 2.5. Copy operation and no-cloning theorem

Most computer programs on classical computers have an assignment of the form  $y = x$ , or  $y = \text{copy}(x)$ , which stores the value in the variable  $x$  in a new location in memory as a variable  $y$ . In scientific computation, this is the foundation of iterative methods, which solve a problem by making progress gradually. For example, classical iterative algorithms for solving linear systems require storing intermediate variables. It is therefore striking that such a basic step is explicitly ruled out by quantum mechanics.

The **no-cloning theorem** is an early result in quantum computation: it forbids a universal quantum copy operation (see also [NC00, Section 12.1]).

**THEOREM 2.38 (No cloning).** *Given a fixed state  $|s\rangle$  (e.g.  $|s\rangle = |0^n\rangle$ ), there is no unitary operator  $U$  that acts as a copy operation, in the sense that for every state  $|x\rangle$ ,*

$$(2.94) \quad U |x\rangle \otimes |s\rangle = |x\rangle \otimes |x\rangle.$$

**PROOF.** Assume such a  $U$  exists. Take two states  $|x_1\rangle, |x_2\rangle$  such that  $0 < |\langle x_1 | x_2 \rangle| < 1$ . Then

$$(2.95) \quad U(|x_1\rangle \otimes |s\rangle) = |x_1\rangle \otimes |x_1\rangle, \quad U(|x_2\rangle \otimes |s\rangle) = |x_2\rangle \otimes |x_2\rangle.$$

Taking the inner product of the two equations and using unitarity,

$$(2.96) \quad \langle x_1 | x_2 \rangle = \langle x_1, s | x_2, s \rangle = \langle x_1, s | U^\dagger U | x_2, s \rangle = \langle x_1, x_1 | x_2, x_2 \rangle = \langle x_1 | x_2 \rangle^2.$$

Hence  $\langle x_1 | x_2 \rangle \in \{0, 1\}$ , contradicting  $0 < |\langle x_1 | x_2 \rangle| < 1$ . □

There are two important special cases in which copying is possible without contradicting Theorem 2.38. The first is that  $|x\rangle$  is not arbitrary: it is a specific state for which we know a preparation procedure, i.e.,  $|x\rangle = U_x |s\rangle$  for a known unitary  $U_x$  and some fixed state  $|s\rangle$ . Then we can prepare a second copy of  $|x\rangle$  via

$$(2.97) \quad (I \otimes U_x) |x\rangle \otimes |s\rangle = |x\rangle \otimes |x\rangle.$$

The second is copying classical information in the computational basis, using the CNOT gate. i.e.,

$$(2.98) \quad \text{CNOT} |x, 0\rangle = |x, x\rangle, \quad x \in \{0, 1\}.$$

The same principle applies to copying classical information from multiple qubits. Fig. 2.4 gives an example of copying the classical information stored in 3 bits.

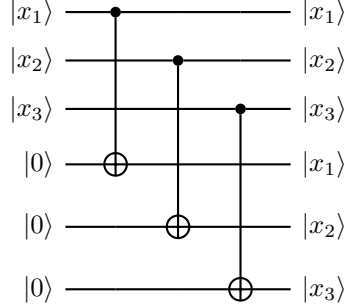


FIGURE 2.4. Copying classical information using multi-qubit CNOT gates.

In general, multi-qubit CNOT operations can be used to perform classical copying in the computational basis. Note that in the circuit model, this can be implemented with a depth-1 circuit, since these CNOT gates act on disjoint sets of qubits.

The copying of classical information is compatible with Theorem 2.38 in the following sense. The proof of Theorem 2.38 uses two non-orthogonal states  $|x_1\rangle, |x_2\rangle$  to obtain a contradiction. However, all states in the computational basis are orthogonal to each other. Therefore, there exist unitaries that copy a specified orthonormal set of states, but a universal quantum copy operation is impossible.

**Example 2.39.** Let us verify that the CNOT gate does not violate the no-cloning theorem, i.e., it cannot be used to copy a general superposition  $|x\rangle = a|0\rangle + b|1\rangle$ . Direct calculation shows

$$(2.99) \quad \text{CNOT } |x\rangle \otimes |0\rangle = a|00\rangle + b|11\rangle \neq |x\rangle \otimes |x\rangle$$

unless  $ab = 0$ . In particular, if  $|x\rangle = |+\rangle$ , then CNOT creates a Bell state.  $\diamond$

Similar to the quantum no-cloning theorem, there does not exist a unitary  $U$  that performs a “deleting” operation which resets an unknown state  $|x\rangle$  to  $|0^n\rangle$ :

$$(2.100) \quad U|0^n\rangle \otimes |x\rangle = |0^n\rangle \otimes |0^n\rangle$$

for all  $|x\rangle$ . Indeed, if  $|x_1\rangle, |x_2\rangle$  are orthogonal, then unitarity implies

$$(2.101) \quad 0 = \langle 0^n, x_1 | 0^n, x_2 \rangle = \langle 0^n, x_1 | U^\dagger U | 0^n, x_2 \rangle = \langle 0^n, 0^n | 0^n, 0^n \rangle = 1,$$

a contradiction.

A more general version of the no-deleting theorem is as follows: given two copies of an arbitrary quantum state, it is impossible to delete one of the copies. Specifically, there is no unitary  $U$  performing the following operation using fixed known states  $|s\rangle, |s'\rangle$ ,

$$(2.102) \quad U|x\rangle|x\rangle|s\rangle = |x\rangle|0^n\rangle|s'\rangle$$

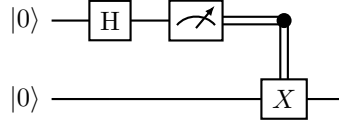
for an arbitrary unknown state  $|x\rangle$ .

**Exercise 2.5.** Prove the version of the no-deleting theorem in Eq. (2.102).

### 2.6. Deferred and implicit measurements

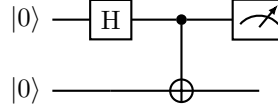
There are two important principles related to quantum measurements: the principle of deferred measurement, and the principle of implicit measurement. At first glance, both principles may seem counterintuitive.

**Example 2.40** (Deferring quantum measurements). Consider the circuit



Here the double line denotes a classical control operation. The outcome is that qubit 0 has probability 1/2 of outputting 0, and qubit 1 is in the state  $|0\rangle$ . Qubit 0 also has probability 1/2 of outputting 1, and qubit 1 is in the state  $|1\rangle$ .

However, we may replace the classical control operation after the measurement by a quantum controlled  $X$  (i.e. CNOT), and measure qubit 0 afterwards:

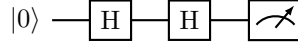


It can be verified that the result is the same. In this sense, CNOT copies the measurement outcome of qubit 0 to qubit 1 in the computational basis.  $\diamond$

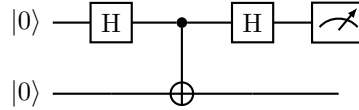
**Example 2.41** (Deferring measurement requires extra qubits). The procedure of deferring quantum measurements using CNOTs is general, and important. Consider the following circuit:



The probability of obtaining 0 or 1 is 1/2. However, if we simply “defer” the measurement to the end by removing the intermediate measurement, we obtain

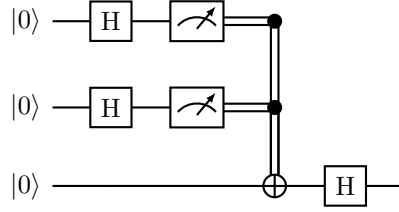


The result of the measurement is deterministically 0! The correct way of deferring the intermediate quantum measurement is to introduce another qubit



Measuring the qubit 0, we obtain 0 or 1 w.p. 1/2, respectively. Hence when deferring quantum measurements, it is necessary to store the intermediate information in extra (ancilla) qubits, even if such information is not used afterwards.  $\diamond$

**Exercise 2.6.** Consider a quantum circuit with three qubits, initially all in state  $|0\rangle$ . The circuit is as follows:

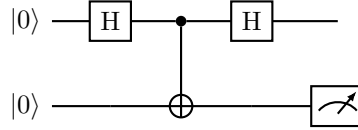


Design a quantum circuit that defers the measurements of the first two qubits to the end, using two additional ancilla qubits to store the intermediate measurement information. After the deferred measurements, describe the final states of all qubits. Ensure the overall effect on the ancilla qubits is the same as if the measurements were performed immediately.

The **principle of deferred measurement** states that in a quantum circuit, measurement operations can be postponed from an intermediate stage to the end of the circuit. This remains true even when a measurement at an intermediate step determines the conditional control of subsequent gates: such classical controls can be replaced by quantum controls. One use of this principle is to simplify quantum circuits and their analysis, by expressing the computation as a unitary circuit (possibly using ancilla qubits) followed by measurements at the end.

The **principle of implicit measurements** states that, for predicting the statistics of the qubits that are measured at the end of a circuit, it is irrelevant whether other qubits are explicitly measured at the end or simply left unmeasured.

**Example 2.42.** Consider the circuit:



Before the measurement, the final state is  $\frac{1}{2}(|00\rangle + |01\rangle) + \frac{1}{2}(|10\rangle - |11\rangle)$ . So measuring qubit 1 yields 0 and 1 with equal probability.

If we measure qubit 0 first, verify that qubit 1 will be in the mixed state

$$(2.103) \quad \rho = \frac{1}{2}|0\rangle\langle 0| + \frac{1}{2}|1\rangle\langle 1|,$$

so if we measure qubit 1 afterwards, we again obtain 0 and 1 with equal probability.  $\diamond$

Why does the principle of implicit measurement hold? Assume the quantum system consists of two subsystems  $A$  and  $B$ . Recall from Eq. (2.74) that a measurement on subsystem  $A$  only depends on the reduced density matrix  $\rho_A$ . Thus it suffices to show that  $\rho_A$  does not depend on whether  $B$  is measured. Let  $\{P_i\}$  be the projectors onto the computational basis of  $B$ , and let the joint state be  $\rho$ . If we measure subsystem  $B$  and discard the outcome, the joint state becomes

$$(2.104) \quad \rho' = \sum_i (I \otimes P_i) \rho (I \otimes P_i).$$

Then

$$(2.105) \quad \rho'_A = \text{Tr}_B[\rho'] = \sum_i \text{Tr}_B[(I \otimes P_i) \rho (I \otimes P_i)] = \sum_i \text{Tr}_B[\rho (I \otimes P_i)] = \text{Tr}_B \left[ \rho \left( I \otimes \sum_i P_i \right) \right] = \text{Tr}_B[\rho] = \rho_A.$$

Therefore, if the qubits in  $A$  are to be measured at the end of the circuit, the measurement statistics do not depend on whether the qubits in  $B$  are measured or not.

### 2.7. Sparse matrix, Majorana, fermionic, and bosonic operators

Sparse matrices are among the most important examples of very large matrices that can be efficiently encoded on quantum computers. They are also closely related to many physical Hamiltonians in practical applications.

**Definition 2.43** ( $s$ -sparse matrix). *A matrix  $A \in \mathbb{C}^{M \times N}$  is called  $s$ -sparse if each row and column of the matrix contains at most  $s$  non-zero entries.*

**Example 2.44.** A diagonal matrix is 1-sparse. Any diagonal matrix  $A \in \mathbb{C}^{2^n \times 2^n}$  can be written as a linear combination of Pauli  $Z$ -operators

$$(2.106) \quad A = \sum_{i_1, \dots, i_n \in \{0,1\}} J_{i_1, \dots, i_n} \sigma_{i_1,1} \cdots \sigma_{i_n,n},$$

where  $\sigma_{s,k}$  is equal to  $Z_k$  if  $s = 1$  and  $I$  if  $s = 0$ . Any permutation matrix  $\Pi$  is 1-sparse. A row and column permutation of a 1-sparse matrix is 1-sparse. Any 1-sparse matrix  $A$  can be written as  $\Pi D$  or  $D\Pi'$ , where  $D$  is a diagonal matrix and  $\Pi, \Pi'$  are permutation matrices. A tridiagonal matrix is 3-sparse. The following matrix

$$(2.107) \quad A = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 1 & 0 & \cdots & 0 \\ \vdots & & & \vdots \\ 1 & 0 & \cdots & 0 \end{pmatrix} = \left( \sum_{i=1}^N e_i \right) e_1^\top \in \mathbb{R}^{N \times N}$$

has only one nonzero entry per row, but it is **not** 1-sparse since the first column has  $N$  nonzero entries.  $\diamond$

**Definition 2.45.** *The maximal absolute value of the entries of  $A \in \mathbb{C}^{M \times N}$ , also called the **max norm**, is defined as:*

$$(2.108) \quad \|A\|_{\max} := \max_{i,j} |A_{ij}|.$$

**Lemma 2.46.** *Let  $A \in \mathbb{C}^{N \times N}$  be  $s$ -sparse. Then*

$$(2.109) \quad \|A\| \leq s \|A\|_{\max}.$$

PROOF. For any row  $i$  of  $A$ , the set of nonzero column indices is denoted by  $\mathcal{C}_i$ . By Cauchy-Schwarz,

$$(2.110) \quad |(Ax)_i|^2 = \left| \sum_{j \in \mathcal{C}_i} A_{ij} x_j \right|^2 \leq \sum_{j \in \mathcal{C}_i} |A_{ij}|^2 \sum_{j \in \mathcal{C}_i} |x_j|^2 \leq s \|A\|_{\max}^2 \sum_{j \in \mathcal{C}_i} |x_j|^2.$$

Then

$$(2.111) \quad \|Ax\|^2 \leq s \|A\|_{\max}^2 \sum_i \sum_{j \in \mathcal{C}_i} |x_j|^2.$$

The condition that  $A$  is  $s$ -sparse implies that for each  $j$ , there are at most  $s$  indices  $i$  such that  $A_{ij} \neq 0$ , i.e.,  $j$  belongs to at most  $s$  sets among  $\{\mathcal{C}_i\}_i$ . Therefore each  $j$  can appear at most  $s$  times in the double sum. This means

$$(2.112) \quad \|Ax\|^2 \leq s^2 \|A\|_{\max}^2 \sum_j |x_j|^2 = s^2 \|A\|_{\max}^2 \|x\|^2.$$

Taking the supremum over  $x \neq 0$  yields  $\|A\| \leq s \|A\|_{\max}$ .  $\square$

The equality in Lemma 2.46 can be reached by considering a matrix  $B$  whose upper left  $s \times s$  block is  $\|A\|_{\max} ee^\top$ , where  $e$  is an all 1 vector of length  $s$ . Direct computation shows that  $\|B\| = s \|A\|_{\max}$ .

A useful lemma is that the product of any 1-sparse matrices is 1-sparse.

**Lemma 2.47.** *Let  $A$  and  $B$  be  $N \times N$  1-sparse matrices. Then  $C = AB$  is also 1-sparse.*

PROOF. Since  $A, B$  are 1-sparse, there exists permutation matrices  $\Pi, \Pi'$  and diagonal matrices  $D, D'$  so that  $A = \Pi D, B = D' \Pi'$ . Therefore

$$(2.113) \quad C = \Pi(DD')\Pi'$$

is a permutation of a diagonal matrix, and is therefore 1-sparse.  $\square$

**Example 2.48.** All Pauli gates in  $\mathcal{P}_n$  are 1-sparse. This can be proved by induction. First, all Pauli matrices  $I, X, Y, Z$  are 1-sparse matrices. Assume all Pauli gates in  $\mathcal{P}_{n-1}$  are 1-sparse, then an element in  $\mathcal{P}_n$  can always be constructed (up to a reordering of qubits) as

$$(2.114) \quad P \otimes P_1, \quad P \in \mathcal{P}_{n-1}, P_1 \in \mathcal{P}_1.$$

This replaces a nonzero entry in  $P$  by a  $2 \times 2$  matrix that is 1-sparse, so the overall matrix is still 1-sparse.  $\diamond$

**Example 2.49** (Majorana operator). For a fermionic system defined on  $n$  modes, the state space  $\mathcal{F} = \otimes_{i=1}^n \mathbb{C}^2 \cong \mathbb{C}^{2^n}$  is called the Fock space. The Majorana fermion operators (or Majorana operators for short) denoted by  $\{\gamma_i\}_{i=1}^{2n}$ , are Hermitian operators in  $L(\mathbb{C}^{2^n})$  satisfying the anticommutation relations:

$$(2.115) \quad \{\gamma_i, \gamma_j\} := \gamma_i \gamma_j + \gamma_j \gamma_i = 2\delta_{ij}, \quad i, j = 1, \dots, 2n.$$

The canonical realization of Majorana operators is through Pauli operators. When  $n = 1$ , we simply have

$$(2.116) \quad \gamma_1 = X, \quad \gamma_2 = Y.$$

For the  $n$  mode system, the Majorana operators can be defined using the **Jordan–Wigner transformation**,

$$(2.117) \quad \gamma_{2j-1} = \left( \prod_{k=1}^{j-1} Z_k \right) X_j, \quad \gamma_{2j} = \left( \prod_{k=1}^{j-1} Z_k \right) Y_j, \quad j = 1, \dots, n.$$

So Majorana operators are also 1-sparse matrices. Furthermore, any product of Majorana operators  $\gamma_{i_1} \cdots \gamma_{i_k}$ ,  $i_1, \dots, i_k \in \{1, \dots, 2n\}$  is 1-sparse.  $\diamond$

**Example 2.50** (Fermionic operator). For a fermionic system defined on  $n$  modes with the Fock space  $\mathcal{F} = \otimes_{i=1}^n \mathbb{C}^2 \cong \mathbb{C}^{2^n}$ , the fermionic creation and annihilation operators, denoted by  $a_i^\dagger$  and  $a_i$  respectively, are operators in  $L(\mathbb{C}^{2^n})$  that satisfy the **canonical anticommutation relations** (CAR):

$$(2.118) \quad \{a_i, a_j^\dagger\} := a_i a_j^\dagger + a_j^\dagger a_i = \delta_{ij}, \quad \{a_i, a_j\} = \{a_i^\dagger, a_j^\dagger\} = 0, \quad i, j = 1, \dots, n.$$

The creation operator  $a_i^\dagger$  adds a fermion to the mode  $i$ , while the annihilation operator  $a_i$  removes a fermion from the mode  $i$ .

For a single mode system,

$$(2.119) \quad a = X^+ = \frac{1}{2}(X + iY) = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, \quad a^\dagger = X^- = \frac{1}{2}(X - iY) = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}.$$

In this convention  $a^\dagger |0\rangle = |1\rangle$ ,  $a^\dagger |1\rangle = 0$ ,  $a |1\rangle = |0\rangle$ ,  $a |0\rangle = 0$ . Here  $|s\rangle$  denotes the state with  $s$  fermions ( $s = 0, 1$ ). The number operator  $\hat{n} = a^\dagger a = \frac{1}{2}(1 - Z)$  satisfies  $\hat{n} |s\rangle = s |s\rangle$ .

For an  $n$ -mode system, the fermionic operators are related to the Majorana operators according to the relation:

$$(2.120) \quad a_i = \frac{1}{2}(\gamma_{2i-1} + i\gamma_{2i}), \quad a_i^\dagger = \frac{1}{2}(\gamma_{2i-1} - i\gamma_{2i}), \quad i = 1, \dots, n,$$

where  $\gamma_{2i-1}$  and  $\gamma_{2i}$  are the Majorana operators associated with the  $i$ -th fermionic mode. Therefore any operator defined using a linear combination of fermionic creation and annihilation operators can be expressed as a linear combination of Majorana operators, and vice versa.

From the Jordan–Wigner transformation,

$$(2.121) \quad a_j = \left( \prod_{k=1}^{j-1} Z_k \right) X_j^+, \quad a_j^\dagger = \left( \prod_{k=1}^{j-1} Z_k \right) X_j^-,$$

with

$$(2.122) \quad X_j^+ = \frac{1}{2}(X_j + iY_j), \quad X_j^- = \frac{1}{2}(X_j - iY_j).$$

Since  $X^\pm$  are 1-sparse matrices,  $a_j^\dagger, a_j, a_j^\dagger a_j, a_j a_j^\dagger$  are also 1-sparse. Furthermore, any product of fermionic operators  $a_{i_1}^\dagger \cdots a_{i_k}^\dagger a_{j_1} \cdots a_{j_l}$  is 1-sparse.  $\diamond$

**Example 2.51** (Bosonic operator). For an  $n$ -mode bosonic systems, the bosonic creation and annihilation operators, denoted by  $b_i^\dagger$  and  $b_i$  respectively, are operators that satisfy the **canonical commutation relations** (CCR):

$$(2.123) \quad [b_i, b_j^\dagger] := b_i b_j^\dagger - b_j^\dagger b_i = \delta_{ij}, \quad [b_i, b_j] = [b_i^\dagger, b_j^\dagger] = 0, \quad i, j = 1, \dots, n.$$

The creation operator  $b_i^\dagger$  adds a boson to the mode  $i$ , while the annihilation operator  $b_i$  removes a boson from the mode  $i$ .

When  $n = 1$ , these operators satisfy

$$(2.124) \quad b |0\rangle = 0, \quad b |s\rangle = \sqrt{s} |s-1\rangle, \quad s = 1, 2, \dots,$$

and

$$(2.125) \quad b^\dagger |s\rangle = \sqrt{s+1} |s+1\rangle, \quad s = 0, 1, 2, \dots$$

Here  $|s\rangle$  denotes a state with  $s$  bosons. We also have

$$(2.126) \quad b^\dagger b |s\rangle = s |s\rangle.$$



In the matrix form, we can write

$$(2.127) \quad b = \begin{pmatrix} 0 & \sqrt{1} & 0 & 0 & \cdots \\ 0 & 0 & \sqrt{2} & 0 & \cdots \\ 0 & 0 & 0 & \sqrt{3} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}, \quad b^\dagger = \begin{pmatrix} 0 & 0 & 0 & 0 & \cdots \\ \sqrt{1} & 0 & 0 & 0 & \cdots \\ 0 & \sqrt{2} & 0 & 0 & \cdots \\ 0 & 0 & \sqrt{3} & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}.$$

These operators are infinite dimensional operators, i.e., operators defined in an infinite dimensional space. They are also 1-sparse. Furthermore,  $\|b\|, \|b^\dagger\| = \infty$ , so unlike any finite dimensional matrices, these operators are unbounded. The physical reason is that a single bosonic mode can accommodate an infinite number of bosons, and the energy of a system with an infinite number of bosons in a single mode is infinity.

Due to the commutation relation, multi-mode bosonic operators can be defined using tensor products:

$$(2.128) \quad b_i = I^{\otimes(i-1)} \otimes b \otimes I^{\otimes(n-i)}, \quad b_i^\dagger = I^{\otimes(i-1)} \otimes b^\dagger \otimes I^{\otimes(n-i)}, \quad i = 1, \dots, n,$$

where the identity operator  $I|s\rangle = |s\rangle$  also acts on an infinite dimensional space.

The precise characterization of the Hilbert space for unbounded operators is beyond the scope of this book. However, if we truncate the state space of each bosonic mode to a finite dimensional space with  $d$  levels, i.e.,  $\mathbb{C}^d$ , the state space of a bosonic system defined on  $n$  modes with  $d$  levels per mode is  $\mathcal{F} = \otimes_{i=1}^n \mathbb{C}^d \cong \mathbb{C}^{d^n}$  and is finite dimensional.

In a single-mode truncated bosonic system,  $b, b^\dagger$  are finite dimensional matrices:

$$(2.129) \quad b = \begin{pmatrix} 0 & \sqrt{1} & 0 & 0 & \cdots & 0 \\ 0 & 0 & \sqrt{2} & 0 & \cdots & 0 \\ 0 & 0 & 0 & \sqrt{3} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \sqrt{d-1} \\ 0 & 0 & 0 & 0 & \cdots & 0 \end{pmatrix}, \quad b^\dagger = \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & 0 \\ \sqrt{1} & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sqrt{2} & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sqrt{3} & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sqrt{d-1} & 0 \end{pmatrix}.$$

These are 1-sparse matrices of size  $d \times d$ . Then the multi-mode operators defined in Eq. (2.128) are 1-sparse matrices. Using Lemma 2.47, the product of any multi-mode bosonic operators  $b_{i_1}^\dagger \cdots b_{i_k}^\dagger b_{j_1} \cdots b_{j_l}$ , where  $b, b^\dagger$  are truncated bosonic creation and annihilation operators defined in Eq. (2.129) are 1-sparse matrices.  $\diamond$

**Exercise 2.7.** Prove that the truncated bosonic creation and annihilation operators defined in Eq. (2.129) satisfy the modified commutation relation

$$(2.130) \quad [b, b^\dagger] = 1 - \frac{d}{(d-1)!} (b^\dagger)^{d-1} (b)^{d-1}.$$

## 2.8. Selected Examples of Hamiltonians in Physics, Chemistry, and Optimization

With the introduction of spin, Majorana, fermionic, and bosonic operators, we can provide several examples of Hamiltonians encountered in applications. Although we will not use all of these examples to illustrate the performance of quantum algorithms, the algorithms in this book can be applied to any of them.

### 2.8.1. Condensed matter physics.

**Example 2.52** (Transverse field Ising model). The Hamiltonian for the one dimensional transverse field Ising model (TFIM) with nearest neighbor interaction of length  $n$  is

$$(2.131) \quad H = - \sum_{i=1}^{n-1} Z_i Z_{i+1} - g \sum_{i=1}^n X_i,$$

where  $g$  is the coupling constant.  $\diamond$

**Example 2.53** (1D Heisenberg model). The Hamiltonian for the 1D Heisenberg model with nearest neighbor interaction is given by

$$(2.132) \quad H = -J \sum_{i=1}^{n-1} \mathbf{S}_i \cdot \mathbf{S}_{i+1}$$

where  $J$  is the interaction strength and  $\mathbf{S}_i$  represents the spin operator at site  $i$ , defined as

$$(2.133) \quad \mathbf{S}_i = \frac{1}{2} \begin{pmatrix} X_i \\ Y_i \\ Z_i \end{pmatrix}.$$

We can decompose this Hamiltonian into three terms, each associated with the  $x$ ,  $y$ , and  $z$  components of the spins:

$$(2.134) \quad H_x = -\frac{J}{4} \sum_{i=1}^{n-1} X_i X_{i+1}, \quad H_y = -\frac{J}{4} \sum_{i=1}^{n-1} Y_i Y_{i+1}, \quad H_z = -\frac{J}{4} \sum_{i=1}^{n-1} Z_i Z_{i+1}.$$

When  $J > 0$  the problem is called ferromagnetic, and when  $J < 0$  it is called anti-ferromagnetic.  $\diamond$

**Example 2.54** (2D Heisenberg model). The Hamiltonian for the 2D Heisenberg model on a square lattice is given by:

$$(2.135) \quad H = -J \sum_{i=1}^{n-1} \sum_{j=1}^{n-1} (\mathbf{S}_{i,j} \cdot \mathbf{S}_{i+1,j} + \mathbf{S}_{i,j} \cdot \mathbf{S}_{i,j+1})$$

We decompose this Hamiltonian into three terms associated with the  $x$ ,  $y$ , and  $z$  components of the spins:

$$(2.136) \quad \begin{aligned} H_x &= -\frac{J}{4} \sum_{i=1}^{n-1} \sum_{j=1}^{n-1} (X_{i,j} X_{i+1,j} + X_{i,j} X_{i,j+1}), \\ H_y &= -\frac{J}{4} \sum_{i=1}^{n-1} \sum_{j=1}^{n-1} (Y_{i,j} Y_{i+1,j} + Y_{i,j} Y_{i,j+1}), \\ H_z &= -\frac{J}{4} \sum_{i=1}^{n-1} \sum_{j=1}^{n-1} (Z_{i,j} Z_{i+1,j} + Z_{i,j} Z_{i,j+1}). \end{aligned}$$

$\diamond$

**Example 2.55** ( $k$ -Local Hamiltonian). A  $k$ -local Hamiltonian is a quantum Hamiltonian where each term acts nontrivially on at most  $k$  qubits. One convenient way to write such a Hamiltonian on  $n$  qubits is as a linear combination of Pauli strings of weight at most  $k$ . For example, one may write

$$(2.137) \quad H = \sum_{S \subseteq [n], |S| \leq k} \sum_{\alpha \in \{0,1,2,3\}^S} J_{S,\alpha} \prod_{i \in S} \sigma_{\alpha(i),i},$$

where  $\sigma_{0,i} = I$ ,  $\sigma_{1,i} = X_i$ ,  $\sigma_{2,i} = Y_i$ , and  $\sigma_{3,i} = Z_i$ .

For example, consider a 2-local Hamiltonian for an  $n$ -qubit system:

$$(2.138) \quad H = \sum_{i < j} J_{ij} \sigma_i \sigma_j,$$

where  $\sigma_i, \sigma_j$  are Pauli operators acting on qubits  $i$  and  $j$ , respectively. Transverse Ising models and Heisenberg models are 2-local Hamiltonians.  $\diamond$

**Example 2.56** (Quadratic fermionic Hamiltonians). Consider the following  $n$ -mode fermionic Hamiltonian

$$(2.139) \quad H = \sum_{k=1}^n \lambda_k c_k^\dagger c_k = \sum_{k=1}^n \frac{\lambda_k}{2} (1 - Z_k),$$

where  $c_k^\dagger$  and  $c_k$  are new fermionic creation and annihilation operators, and  $\lambda_k$  are real eigenvalues representing the energy levels of the system. The Hamiltonian  $H$  is a linear combination of Pauli  $Z$  operators and is thus a diagonal matrix.

Now, consider a general quadratic fermionic Hamiltonian of the form:

$$(2.140) \quad H = \sum_{i,j=1}^n A_{ij} a_i^\dagger a_j,$$

where  $A$  is a Hermitian matrix. Since  $A$  is Hermitian, we can diagonalize it using a unitary transformation  $U$  such that:

$$(2.141) \quad U^\dagger A U = \Lambda,$$

where  $\Lambda$  is a diagonal matrix containing the eigenvalues  $\lambda_k$ . Then define

$$(2.142) \quad c_k = \sum_{i=1}^n (U^\dagger)_{ki} a_i, \quad c_k^\dagger = \sum_{i=1}^n a_i^\dagger U_{ik}, \quad k = 1, \dots, n.$$

Direct calculation shows that the new set of creation and annihilation operators  $\{c_k^\dagger, c_k\}$  satisfy the canonical anticommutation relation. Substituting these transformations into the Hamiltonian,

$$(2.143) \quad H = \sum_{i,j,k} U_{ik} \lambda_k (U^\dagger)_{kj} a_i^\dagger a_j = \sum_{k=1}^n \lambda_k c_k^\dagger c_k,$$

we have transformed  $H$  into a diagonal Hamiltonian.  $\diamond$

**Example 2.57** (1D spinless Hubbard model). The Hamiltonian for the 1D spinless Hubbard model with nearest-neighbor interaction is given by:

$$(2.144) \quad H = -t \sum_{i=1}^{n-1} (a_i^\dagger a_{i+1} + a_{i+1}^\dagger a_i) + U \sum_{i=1}^{n-1} n_i n_{i+1},$$

where  $t$  is the hopping parameter, representing the kinetic energy term, and  $U$  is the nearest-neighbor interaction strength. The operators  $a_i^\dagger$  and  $a_i$  are the fermionic creation and annihilation operators at site  $i$ , respectively, and  $n_i = a_i^\dagger a_i$  is the number operator at site  $i$ . When  $U = 0$ , the Hamiltonian is a quadratic in the fermionic operators and can be turned into a diagonalized form. When  $U \neq 0$ , the Hamiltonian is no longer quadratic and cannot be turned into a diagonalized Hamiltonian using the same strategy.  $\diamond$

**Example 2.58** (Uniform electron gas in a plane wave basis). In a plane wave basis, the Hamiltonian for a box of uniform electron gas can be expressed in second quantization as follows:

$$(2.145) \quad H = \sum_{\mathbf{k}} \epsilon_{\mathbf{k}} c_{\mathbf{k}}^\dagger c_{\mathbf{k}} + \frac{1}{2} \sum_{\mathbf{k}_1, \mathbf{k}_2, \mathbf{q}} V(\mathbf{q}) c_{\mathbf{k}_1 + \mathbf{q}}^\dagger c_{\mathbf{k}_2 - \mathbf{q}}^\dagger c_{\mathbf{k}_2} c_{\mathbf{k}_1},$$

where  $c_{\mathbf{k}}^\dagger$  and  $c_{\mathbf{k}}$  are fermionic creation and annihilation operators for an electron with wave vector  $\mathbf{k} \in \mathbb{R}^3$ ,  $\epsilon_{\mathbf{k}} = |\mathbf{k}|^2/2$  is the kinetic energy. The interaction potential  $V(\mathbf{q}) = 4\pi/\mathbf{q}^2$  in a plane wave basis is the Fourier transform of the Coulomb potential.  $\diamond$

**Example 2.59** (Harmonic oscillator). The Hamiltonian for a quantum harmonic oscillator in the first quantization (i.e., real space representation) is given by

$$(2.146) \quad H = \frac{p^2 + x^2}{2},$$

where  $p = -i\partial_x$  is the momentum operator and  $x$  is the position operator. Define

$$(2.147) \quad b = \frac{1}{\sqrt{2}}(x + ip), \quad b^\dagger = \frac{1}{\sqrt{2}}(x - ip),$$

then  $b, b^\dagger$  satisfy the canonical commutation relation  $[b, b^\dagger] = 1$ . Furthermore, the Hamiltonian takes the form

$$(2.148) \quad H = b^\dagger b + \frac{1}{2}.$$

If we truncate the bosonic mode to include  $d$  levels, the state space is  $\mathcal{F} = \mathbb{C}^d$ , and  $H$  is a diagonal matrix of size  $d \times d$ .  $\diamond$

### 2.8.2. Quantum chemistry.

**Example 2.60** (Quantum chemistry in first quantization). In first quantization, the Hamiltonian for a many-electron system is given in terms of the coordinates and momenta of the electrons. The non-relativistic electronic Hamiltonian for a molecule in atomic units can be expressed as:

$$(2.149) \quad H = - \sum_{i=1}^N \frac{\nabla_i^2}{2} - \sum_{i=1}^N \sum_{A=1}^M \frac{Z_A}{|\mathbf{r}_i - \mathbf{R}_A|} + \sum_{i < j} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} + \sum_{A < B} \frac{Z_A Z_B}{|\mathbf{R}_A - \mathbf{R}_B|},$$

where  $N$  is the number of electrons,  $M$  is the number of nuclei,  $\mathbf{r}_i$  and  $\mathbf{R}_A$  are the positions of the  $i$ -th electron and the  $A$ -th nucleus, respectively,  $Z_A$  is the atomic number of the  $A$ -th nucleus. This is an unbounded operator.  $\diamond$

**Example 2.61** (Quantum chemistry in second quantization). In quantum chemistry, the electronic structure of molecules can be described using the formalism of second quantization with  $n$  molecular orbitals. The state space  $\mathcal{F} = \otimes_{i=1}^n \mathbb{C}^2$  is finite dimensional. The use of second quantization allows

for a compact and efficient representation of the Hamiltonian and facilitates the expression of the Hamiltonian on quantum computers via the Jordan–Wigner transformation. The Hamiltonian of a many-electron system in second quantization is given by

$$(2.150) \quad H = \sum_{p,q=1}^n h_{pq} a_p^\dagger a_q + \frac{1}{2} \sum_{p,q,r,s=1}^n V_{pqrs} a_p^\dagger a_q^\dagger a_r a_s,$$

where  $a_p^\dagger$  and  $a_q$  are fermionic creation and annihilation operators, respectively. The creation operator  $a_p^\dagger$  adds an electron to the molecular orbital  $p$ , and the annihilation operator  $a_q$  removes an electron from the molecular orbital  $q$ . For simplicity we only consider the spatial part of the orbital and omit the spin part. The indices  $p, q, r$ , and  $s$  label the molecular orbitals,  $h_{pq}$  are the one-electron integrals, and  $V_{pqrs}$  are the two-electron integrals.

The one-electron integrals  $h_{pq}$  are given by

$$(2.151) \quad h_{pq} = \int \psi_p^*(\mathbf{r}) \left( -\frac{\nabla^2}{2} + V_{\text{ext}}(\mathbf{r}) \right) \psi_q(\mathbf{r}) d\mathbf{r},$$

where  $\psi_p(\mathbf{r})$  is the spatial part of the molecular orbital and  $V_{\text{ext}}(\mathbf{r}) = -\sum_{A=1}^M \frac{Z_A}{|\mathbf{r}-\mathbf{R}_A|}$  is the external potential due to the nuclei. The two-electron integrals  $V_{pqrs}$  are given by

$$(2.152) \quad V_{pqrs} = \int \int \psi_p^*(\mathbf{r}_1) \psi_q^*(\mathbf{r}_2) \frac{1}{|\mathbf{r}_1 - \mathbf{r}_2|} \psi_r(\mathbf{r}_2) \psi_s(\mathbf{r}_1) d\mathbf{r}_1 d\mathbf{r}_2.$$

The nuclei-nuclei interaction is a constant and is dropped for simplicity.  $\diamond$

**Example 2.62** (PPP Model). The Pariser-Parr-Pople (PPP) model is used in quantum chemistry to describe the  $\pi$ -electron systems in conjugated organic molecules. The Hamiltonian for the PPP model can be written as

$$(2.153) \quad H = \sum_{p,q=1}^n h_{pq} a_p^\dagger a_q + \frac{1}{2} \sum_{p,q=1}^n V_{pq} n_p n_q,$$

where  $h_{pq}$  are hopping integral elements,  $V_{pq}$  are Coulomb interaction elements,  $a_p^\dagger$  and  $a_p$  are the fermionic creation and annihilation operators at site  $p$ , and  $n_p = a_p^\dagger a_p$  is the number operator. The Hubbard model is a special case of the PPP model with short ranged hopping and Coulomb interaction elements. Compared to the full chemistry Hamiltonian in second quantization, the two-body interaction coefficients  $V_{pq}$  have only  $\mathcal{O}(n^2)$  entries but can still represent long range interactions.  $\diamond$

### 2.8.3. Quantum field theory.

**Example 2.63** (Schwinger Model in 1D). The Schwinger model describes quantum electrodynamics in  $1+1$  dimensions. The state space for the Schwinger model is the tensor product of two spaces: a tensor product of  $n+1$  fermionic spaces and a product of  $n$  gauge field spaces. The total Fock space is given by

$$(2.154) \quad \mathcal{F} = \left( \bigotimes_{i=1}^{n+1} \mathbb{C}^2 \right) \otimes \left( \bigotimes_{j=1}^n \mathbb{C}^d \right),$$

where  $d = 2L + 1$  is the number of levels the gauge field can take. There are two operators that we need to define that act on the gauge field space. The first is  $E_j^2$ , which is a diagonal operator that

counts the energy stored in the gauge field with index  $j \in \{1, \dots, n\}$ . The second is  $U_j$ , which adds one to the value stored in the gauge field register and is analogous to a bosonic creation operator. The action of these operators is given formally below:

$$(2.155) \quad E_j^2 = \sum_{\varepsilon=-L}^L \varepsilon^2 |\varepsilon\rangle_j \langle \varepsilon|_j, \quad U_j = \sum_{\varepsilon=-L}^L |\varepsilon+1\rangle_j \langle \varepsilon|_j, \quad U_j^\dagger = \sum_{\varepsilon=-L}^L |\varepsilon-1\rangle_j \langle \varepsilon|_j.$$

Here we assume for  $U_j$  and its adjoint that the gauge field satisfies periodic boundary conditions at the cutoff located at  $\varepsilon = \pm L$ .

The Hamiltonian for the Schwinger model is given by:

$$(2.156) \quad H = \sum_{j=1}^n E_j^2 \otimes I_2^{\otimes(n+1)} + \nu \sum_{j=1}^n \left[ U_j \otimes a_j^\dagger a_{j+1} - U_j^\dagger \otimes a_j a_{j+1}^\dagger \right] + \mu \sum_{j=1}^n (-1)^j I_d^{\otimes n} \otimes a_j^\dagger a_j,$$

where  $a_i$  and  $a_i^\dagger$  are the fermionic annihilation and creation operators at site  $i$ , and  $I_m$  denotes the identity operator of dimension  $m$ . The parameters  $\mu, \nu$  are related to parameters such as the lattice spacing.  $\diamond$

**Example 2.64** (Quadratic Majorana operators). From the Jordan–Wigner transformation in Eq. (2.117), and use the fact that  $XY = iZ$ , we find that

$$(2.157) \quad H = -i \sum_{k=1}^n \lambda_k \gamma_{2k-1} \gamma_{2k} = \sum_{k=1}^n \lambda_k Z_k, \quad \lambda_k \in \mathbb{R}$$

is a diagonal Hamiltonian.

Consider a quadratic Hamiltonian of the form:

$$(2.158) \quad H = -i \sum_{1 \leq p < q \leq 2n} A_{pq} \zeta_p \zeta_q = -\frac{i}{2} \sum_{p,q=1}^{2n} A_{pq} \zeta_p \zeta_q,$$

where  $A$  is a real antisymmetric matrix, and  $\{\zeta_p\}_{p=1}^{2n}$  is a set of Majorana operators. There exists an orthogonal matrix  $O$  such that:

$$(2.159) \quad O^\top A O = \bigoplus_{k=1}^n \begin{pmatrix} 0 & \lambda_k \\ -\lambda_k & 0 \end{pmatrix} =: \Lambda$$

where  $\lambda_k$  are the singular values of  $A$ . Now define a set of transformed Majorana operators

$$(2.160) \quad \gamma_j = \sum_p \zeta_p O_{pj} = \sum_p (O^\top)_{jp} \zeta_p, \quad j = 1, \dots, 2n,$$

then we still have

$$(2.161) \quad \{\gamma_j, \gamma_{j'}\} = 2\delta_{j,j'}.$$

The transformed Hamiltonian takes a diagonal form

$$(2.162) \quad H = -\frac{i}{2} \sum_{1 \leq j, j' \leq 2n} \gamma_j (\Lambda)_{jj'} \gamma_{j'} = -i \sum_{k=1}^n \lambda_k \gamma_{2k-1} \gamma_{2k}.$$

The quadratic fermionic Hamiltonian in Example 2.56 is a special case of this example.  $\diamond$

**Example 2.65** (SYK Model). The Sachdev-Ye-Kitaev (SYK) model is a quantum mechanical model of  $n$  Majorana fermions with random all-to-all interactions. The Hamiltonian for the SYK model is given by

$$(2.163) \quad H = \sum_{1 \leq i < j < k < l \leq 2n} J_{ijkl} \gamma_i \gamma_j \gamma_k \gamma_l,$$

where  $\gamma_i$  are the Majorana fermion operators, and  $J_{ijkl}$  are random coupling constants, typically drawn from a Gaussian distribution. The SYK model is of particular interest due to its connections to quantum chaos, holography, and black hole physics.  $\diamond$

#### 2.8.4. Optimization.

**Example 2.66** ( $k$ -SAT problem). Classical optimization problems, such as the  $k$ -SAT problem, can be represented using a Hamiltonian. The  $k$ -SAT problem is a type of Boolean satisfiability problem where each clause contains exactly  $k$  literals. The goal is to find an assignment to the Boolean variables that satisfies all the clauses. The most famous examples are 2-SAT (classically easy), and 3-SAT (NP-complete).

Consider a  $k$ -SAT problem with  $n$  Boolean variables  $x_1, x_2, \dots, x_n$  and  $m$  clauses  $C_1, C_2, \dots, C_m$ . Each clause  $C_i$  is a disjunction of exactly  $k$  literals.

We can construct a Hamiltonian  $H$  such that its ground state corresponds to the solution of the  $k$ -SAT problem. The Hamiltonian for the  $k$ -SAT problem can be written as:

$$(2.164) \quad H = \sum_{i=1}^m H_{C_i},$$

where  $H_{C_i}$  is the Hamiltonian for the  $i$ -th clause. Each clause Hamiltonian  $H_{C_i}$  is designed to be zero if the clause is satisfied and positive otherwise. For clauses involving single literals, such as  $C_k = (x_p)$  or  $C_l = (\bar{x}_q)$ , the Hamiltonians  $H_{C_k}$  and  $H_{C_l}$  are:

$$(2.165) \quad H_{C_k} = \frac{1}{2} (1 + Z_p), \quad H_{C_l} = \frac{1}{2} (1 - Z_q).$$

For a clause  $C_i = (x_p \vee \bar{x}_q)$ , the corresponding Hamiltonian  $H_{C_i}$  can be written using the product

$$(2.166) \quad H_{C_i} = \frac{1}{4} (1 + Z_p) (1 - Z_q).$$

For a general clause  $C_i = (l_1 \vee l_2 \vee \dots \vee l_k)$ , where  $l_j$  represents either  $x_{p_j}$  or  $\bar{x}_{p_j}$ , the corresponding Hamiltonian  $H_{C_i}$  can be written using the Pauli-Z operator  $Z$ :

$$(2.167) \quad H_{C_i} = \prod_{j=1}^k \frac{1 + z_j Z_{p_j}}{2},$$

where  $z_j = +1$  if  $l_j = x_{p_j}$  and  $z_j = -1$  if  $l_j = \bar{x}_{p_j}$ . The Hamiltonian  $H$  is diagonal and positive semidefinite. If the smallest eigenvalue (called the ground state energy) of  $H$  is 0, then the associated eigenvector (called the ground state, which may not be unique) corresponds to the Boolean variable assignment that satisfies all the clauses of the  $k$ -SAT problem.  $\diamond$

**Example 2.67** (MAX-CUT problem). The MAX-CUT problem is a well-known combinatorial optimization problem. Given a graph  $G = (V, E)$  with a set of vertices  $V$  and a set of edges  $E$ , the

goal is to partition the vertices into two subsets such that the number of edges between the subsets is maximized. Assume the graph has  $n$  vertices, and the Hamiltonian for the MAX-CUT problem can be written as:

$$(2.168) \quad H = - \sum_{(i,j) \in E} \frac{1}{2} (1 - Z_i Z_j).$$

Each term  $-\frac{1}{2}(1 - Z_i Z_j)$  equals  $-1$  if vertices  $i$  and  $j$  are in different subsets and  $0$  if they are in the same subset. Therefore, minimizing  $H$  is equivalent to maximizing the number of edges that are cut by the partition.  $\diamond$



Part II

Foundation



## CHAPTER 3

# Probability, quantum channel, and distances

We begin by reviewing basic concepts in classical probability theory, which provides intuition for how errors propagate in randomized processes. We then introduce quantum channels as the general framework for quantum dynamics. Unlike ideal quantum circuits which are unitary, real-world quantum processes often involve noise and decoherence. Quantum channels allow us to model these effects, as well as measurements and interactions with the environment. We explain the requirements (specifically, the concept of complete positivity) for a map to be a valid quantum channel and describe standard representations such as the Kraus and Stinespring forms.

With this framework in place, we introduce distance measures for quantum states. For pure states, we use norms that account for the global phase. For mixed states, we introduce the trace distance and fidelity. These two measures are complementary: trace distance relates to the distinguishability of states via measurement, while fidelity captures their overlap and behaves well under quantum operations.

Finally, we discuss how to compare quantum channels. This requires norms that are stable even when the channels act on part of an entangled system. This leads us to the diamond norm, which is the standard metric for quantifying the error of quantum operations.

### 3.1. Basic notions in probability theory

Probability theory is a subject that carries nearly as many profound surprises as quantum theory itself. In this section, we introduce some basic concepts in probability theory, focusing on finite-dimensional spaces. In quantum computing, the probability distributions associated with an  $n$ -qubit system reside in  $2^n$ -dimensional spaces.

**Definition 3.1.** *Let  $\Sigma$  be a finite set called a sample space, where each element of  $\Sigma$  is called an event. A **probability distribution** is a function  $\mathbb{P} : \Sigma \rightarrow [0, 1]$ , which can be represented as a vector in a Euclidean space, and satisfies  $\sum_{s \in \Sigma} \mathbb{P}(s) = 1$ .*

Let  $\Sigma_A$  and  $\Sigma_B$  be sample spaces and let  $\mathbb{P}_A$  and  $\mathbb{P}_B$  be probability distributions on the two sample spaces. These distributions are said to be independent if the joint distribution,  $\mathbb{P}_{AB}$  on the set  $\Sigma_A \times \Sigma_B$  obeys  $\mathbb{P}_{AB} = \mathbb{P}_A \otimes \mathbb{P}_B$ . The expectation value (or average value) of a function mapping  $f : \Sigma \mapsto \mathbb{C}$  is defined to be  $\mathbb{E}(f) := \sum_{s \in \Sigma} f(s) \mathbb{P}(s) = \langle f, \mathbb{P} \rangle$ .

**Example 3.2.** As an example, let us consider rolling a four-sided die. Here the random variable is the outcome of the experiment; the sample space is  $\{1, 2, 3, 4\}$  and the probability distribution (for a fair die) is  $1/4$  for each of these outcomes. The random variable,  $x$ , in this case corresponds to the result of the die.

In the event that we wanted to find the probability that the sample is a prime number, we could redefine the sample space and the underlying distribution but it is easier to use the indicator-function property of the distribution to see that

$$(3.1) \quad \mathbb{P}(x \in \{2, 3\}) = \mathbb{E}(\mathbf{1}_{\{2,3\}}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}.$$

In general, this approach is often the easiest way to compute a probability because it constructs an indicator function which projects onto the fraction of the sample space that we want to measure. Note this is also true in quantum theory wherein the probability of measuring a mixed state,  $\rho$ , to be a pure state  $|\psi\rangle$  is

$$(3.2) \quad \mathbb{P}(|\psi\rangle) = \text{Tr}(|\psi\rangle\langle\psi|\rho) = \langle\psi|\rho|\psi\rangle.$$

Here the projector  $|\psi\rangle\langle\psi|$  plays the same role as the indicator function used above, and further illustrates the close ties between probability theory and quantum theory.  $\diamond$

Similar to the amplitude of the wave function in quantum theory, there is not a single unifying interpretation of probability. For this reason we recommend that the reader be well versed in both interpretations as each can convey useful intuitions.

The following bound, known as the union bound, is very useful for estimating probabilities of events. We provide it as well as its proof as an elementary example of probability theory.

**THEOREM 3.3 (Union Bound).** *Let  $\Sigma$  be a sample space and let  $A, B \subseteq \Sigma$  and let  $\mathbb{P}$  be a probability distribution on  $\Sigma$ . We then have*

$$(3.3) \quad \mathbb{P}(A \cup B) = \mathbb{E}(\mathbf{1}_A + \mathbf{1}_B - \mathbf{1}_A \mathbf{1}_B) \leq \mathbb{P}(A) + \mathbb{P}(B).$$

**PROOF.** Intuitively, by looking at a Venn diagram for events  $A$  and  $B$  it is clear that the region  $A \cup B$  contains region  $A$  and region  $B$  but also may include region  $A \cap B$ . Thus the upper bound given above overcounts the probability in the intersection and therefore it is an upper bound. Formally, we use linearity of expectation:

$$(3.4) \quad \mathbb{E}(\mathbf{1}_A + \mathbf{1}_B - \mathbf{1}_A \mathbf{1}_B) = \mathbb{E}(\mathbf{1}_A) + \mathbb{E}(\mathbf{1}_B) - \mathbb{E}(\mathbf{1}_A \mathbf{1}_B).$$

Next,  $\mathbb{E}(\mathbf{1}_A \mathbf{1}_B) = \sum_{s \in \Sigma} \mathbb{P}(s)(\mathbf{1}_A(s)\mathbf{1}_B(s)) \geq 0$ , and  $\mathbb{E}(\mathbf{1}_A) = \mathbb{P}(A)$ ,  $\mathbb{E}(\mathbf{1}_B) = \mathbb{P}(B)$ . Combining these gives the claim.  $\square$

**Example 3.4 (Failure Propagation Bound).** Consider the following problem: you have a quantum algorithm that succeeds with probability  $1 - \delta$  and fails with probability  $\delta$ . Suppose we run the algorithm independently  $N$  times; determine a value of  $\delta$  that guarantees the probability of at least one failure is at most  $1/3$ . This problem appears ubiquitously in quantum computing in problems such as phase estimation or quantum error correction where the probability of failure needs to be considered and extra computational resources are needed to suppress them.

The  $N$  events each have a probability of  $\delta$  assigned to them and so we expect that the total probability of at least one error happening will be from the union bound  $N\delta$ . We can validate this inductively. For the base case we see trivially that the claim holds for  $N = 1$ . For the induction step, let us assume that the probability of at least one error occurring in the first  $N - 1$  steps is at most  $(N - 1)\delta$ . From the union bound the probability of failing in the next sample is  $\delta$  and thus the total failure probability is at most  $(N - 1)\delta + \delta = N\delta$ . Thus if we want to see a failure probability of  $1/3$  it suffices to take

$$(3.5) \quad \delta \leq \frac{1}{3N}.$$

This example shows that worst case scenario that the failure probability for our algorithm grows linearly. This actually might seem strange to the reader since the error probability compounds exponentially in practice; however, linear growth of error is actually in this context worst than exponential because for large enough  $N$  the union bound will be greater than 1 whereas the exponential upper bound is always less than 1. In this context, surprisingly, linear growth is worse than exponential but nonetheless the simplicity and generality of union bounds often provide good enough bounds that are easy to manipulate.  $\diamond$

The natural operations on probability distributions are stochastic transformations, which can be represented as transition matrices. We define these transformations below.

**Definition 3.5.** Let  $\Sigma$  be a sample space of size  $N$  and let  $p \in \mathbb{R}^N$  be the column vector representation of a probability distribution. A valid transformation on the state space of the register  $X$  to itself has a matrix representation  $P : \mathbb{R}^N \rightarrow \mathbb{R}^N$ , which maps  $p$  to  $Pp$ . The matrix  $P$  is called a **transition matrix** and satisfies

- (1)  $P_{ij} \geq 0, \quad \forall i, j \in [N],$
- (2)  $\sum_{i \in [N]} P_{ij} = 1, \quad \forall j \in [N].$

**Remark 3.6.** In classical probability theory, the probability distribution is often written as a row vector. Then the transition matrix is applied from the right as  $pP$ , and the transition matrix needs to be **right stochastic** or **row stochastic**, i.e.,  $\sum_{j \in [N]} P_{ij} = 1$  for all  $i \in [N]$ . Given a probability distribution  $p \in \mathbb{R}^N$ , a natural quantum state encoding the distribution  $p$  (also called a coherent version of  $p$ ) is

$$(3.6) \quad |\sqrt{p}\rangle = \sum_i \sqrt{p_i} |i\rangle.$$

This is a normalized state. It is thus more natural to view  $p$  as a column vector so that the usual rule of applying an operator to a state vector applies. A matrix satisfying the properties in Definition 3.5 is also called **left stochastic** or **column stochastic**. Any  $j$ -th column of  $P$ , denoted by  $P_{:,j}$ , is a probability distribution. If  $P$  is both left and right stochastic, then it is called a **doubly stochastic** matrix.  $\diamond$

**Example 3.7.** Let us consider how we would represent an AND gate in this language. The AND gate has the property that for any  $x, y \in \{0, 1\}$ ,  $\text{AND}(x, y) = xy$ . This operation is an example of an irreversible operation, meaning that it cannot be inverted from the outputs to find the inputs. In this case the natural vector space for probability distributions for two bits can be represented as a probability vector in  $\mathbb{R}^2 \otimes \mathbb{R}^2$ . As we are using square matrices to represent these transformations we will take  $\text{AND}(e_x \otimes e_y) = e_0 \otimes e_{xy}$  for computational basis vectors  $e_0, e_1$  and  $x, y \in \{0, 1\}$ . Specifically then we have that the gate can be represented as a stochastic matrix  $P_{\text{AND}}$

$$(3.7) \quad P_{\text{AND}} = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}.$$

We see that the matrix representation is stochastic, but not doubly stochastic.

If we consider taking two distributions for our bits  $p_x = [a, 1 - a]^\top$  and  $p_y = [b, 1 - b]^\top$  for  $a, b \in [0, 1]$  then we can see that the distribution that we get from applying the AND operation to

the distribution on the bits is

$$(3.8) \quad P_{\text{AND}}(p_x \otimes p_y) = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} ab \\ a(1-b) \\ b(1-a) \\ (1-a)(1-b) \end{pmatrix} = \begin{bmatrix} (a+b) - ab \\ 1 - (a+b) + ab \\ 0 \\ 0 \end{bmatrix}.$$

This output distribution makes intuitive sense. The AND output is 1 only if both inputs are 1, which occurs with probability  $(1-a)(1-b)$ , corresponding to the second entry above. Equivalently, the probability that the AND output is 0 is the probability that at least one input is 0, namely  $a + b - ab$ , corresponding to the first entry.  $\diamond$

### 3.2. Quantum Channels

The concept of a quantum channel generalizes both the unitary evolution of isolated quantum systems, as governed by the Schrödinger equation, and the stochastic evolution of classical probability distributions. It provides a unified framework for describing the most general physically permissible evolution of quantum states, encompassing coherent dynamics (e.g., unitary transformations) and incoherent processes such as measurement, decoherence, and interactions with an environment.

We begin by defining the mathematical objects under consideration. A **superoperator** is a linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$ . We denote the action of  $\mathcal{Q}$  on an operator  $A \in L(\mathbb{C}^N)$  by  $\mathcal{Q}[A]$  or  $\mathcal{Q}(A)$ .

Given two superoperators  $\mathcal{Q}_1 : L(\mathbb{C}^{N_1}) \rightarrow L(\mathbb{C}^{M_1})$  and  $\mathcal{Q}_2 : L(\mathbb{C}^{N_2}) \rightarrow L(\mathbb{C}^{M_2})$ , their **tensor product**  $\mathcal{Q}_1 \otimes \mathcal{Q}_2$  is the unique linear map  $L(\mathbb{C}^{N_1} \otimes \mathbb{C}^{N_2}) \rightarrow L(\mathbb{C}^{M_1} \otimes \mathbb{C}^{M_2})$  satisfying

$$(3.9) \quad (\mathcal{Q}_1 \otimes \mathcal{Q}_2)[A_1 \otimes A_2] = \mathcal{Q}_1[A_1] \otimes \mathcal{Q}_2[A_2]$$

for all  $A_1 \in L(\mathbb{C}^{N_1})$  and  $A_2 \in L(\mathbb{C}^{N_2})$ . This definition extends to all operators by linearity.

Just as a unitary transformation maps a state vector to another state vector while preserving its norm, a **quantum channel** is a superoperator intended to map a density operator to another density operator. A fundamental example is the **identity channel**  $\mathcal{I} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^N)$ , defined by  $\mathcal{I}[A] = A$  for any  $A \in L(\mathbb{C}^N)$ .

**Example 3.8.** The action of the tensor product of superoperators is particularly important when analyzing local operations on composite systems. Let  $\mathcal{I}_K : L(\mathbb{C}^K) \rightarrow L(\mathbb{C}^K)$  be the identity map and  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  be a linear map. Consider an operator  $A \in L(\mathbb{C}^K \otimes \mathbb{C}^N)$ . We can represent  $A$  in block form with respect to an orthonormal basis  $\{|i\rangle\}$  of  $\mathbb{C}^K$ :

$$(3.10) \quad A = \sum_{i,j \in [K]} |i\rangle\langle j| \otimes A_{ij}, \quad A_{ij} \in L(\mathbb{C}^N).$$

The action of  $\mathcal{I}_K \otimes \mathcal{Q}$  is given by applying  $\mathcal{Q}$  to each block:

$$(3.11) \quad (\mathcal{I}_K \otimes \mathcal{Q})[A] = \sum_{i,j \in [K]} |i\rangle\langle j| \otimes \mathcal{Q}[A_{ij}].$$

For instance, if  $K = 2$ , the matrix representation is

$$(3.12) \quad (\mathcal{I}_2 \otimes \mathcal{Q}) \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} = \begin{pmatrix} \mathcal{Q}[A_{00}] & \mathcal{Q}[A_{01}] \\ \mathcal{Q}[A_{10}] & \mathcal{Q}[A_{11}] \end{pmatrix} \in L(\mathbb{C}^2 \otimes \mathbb{C}^M).$$

$\diamond$

To ensure that a superoperator maps density operators (which are positive semidefinite and have unit trace) to density operators, it must satisfy certain constraints.

**Definition 3.9.** A linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  is called: **positive** if  $\mathcal{Q}[A]$  is positive semidefinite for every positive semidefinite  $A \in L(\mathbb{C}^N)$ .  $\mathcal{Q}$  is called **trace preserving (TP)** if  $\text{Tr}(\mathcal{Q}[A]) = \text{Tr}(A)$  for every  $A \in L(\mathbb{C}^N)$ .

While it might seem sufficient to define a quantum channel simply as a positive, trace-preserving map, the structure of quantum mechanics demands a stronger condition. Quantum systems often exist as subsystems of larger, composite systems. If  $\mathcal{Q}$  describes the evolution of a system  $S$ , and  $S$  is potentially entangled with an ancillary system  $A$ , the evolution of the joint system is described by  $\mathcal{I}_A \otimes \mathcal{Q}$ . For this joint evolution to be physically valid,  $\mathcal{I}_A \otimes \mathcal{Q}$  must also map density operators to density operators, meaning it must be a positive map, regardless of the dimension of the ancilla  $A$ . This requirement leads to the concept of complete positivity.

**Definition 3.10.** A linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  is **completely positive (CP)** if for all integers  $K \geq 1$ , the map  $\mathcal{I}_K \otimes \mathcal{Q} : L(\mathbb{C}^K \otimes \mathbb{C}^N) \rightarrow L(\mathbb{C}^K \otimes \mathbb{C}^M)$  is positive.

It is worth noting that the ordering of the tensor product in the definition is immaterial. One could equivalently require that  $\mathcal{Q} \otimes \mathcal{I}_K$  be positive for all  $K$ . Physically, this reflects the fact that the labeling of the ancillary system is arbitrary. Mathematically, the maps  $\mathcal{I} \otimes \mathcal{Q}$  and  $\mathcal{Q} \otimes \mathcal{I}$  are related via the SWAP operator (the isomorphism that exchanges the tensor factors). Specifically, they are unitarily equivalent:

$$(3.13) \quad \mathcal{Q} \otimes \mathcal{I} = \mathcal{U}_{\text{SWAP}} \circ (\mathcal{I} \otimes \mathcal{Q}) \circ \mathcal{U}_{\text{SWAP}}^{-1},$$

where the superoperator  $\mathcal{U}_{\text{SWAP}}$  acts as  $\mathcal{U}_{\text{SWAP}}[X] = \text{SWAP} \cdot X \cdot \text{SWAP}^\dagger$ . Since  $X \succeq 0$  if and only if  $UXU^\dagger \succeq 0$  for any unitary  $U$ , it follows that  $\mathcal{I} \otimes \mathcal{Q}$  is positive if and only if  $\mathcal{Q} \otimes \mathcal{I}$  is positive.

While positivity ensures that the channel acts correctly on the system itself, complete positivity is strictly stronger, ensuring correct action even when the system is entangled with an ancilla.

**Example 3.11** (Positive map that is not completely positive). Consider the transpose map  $\mathcal{T} : L(\mathbb{C}^2) \rightarrow L(\mathbb{C}^2)$ , defined by  $\mathcal{T}[A] = A^\top$  with respect to the computational basis. If  $A$  is positive, its eigenvalues are non-negative. Since  $A$  and  $A^\top$  share the same spectrum,  $A^\top$  is also positive. Thus,  $\mathcal{T}$  is a positive map.

However,  $\mathcal{T}$  is not completely positive. To illustrate this, consider a two-qubit system in the maximally entangled Bell state  $|\psi\rangle = \frac{1}{\sqrt{2}}(|00\rangle + |11\rangle)$ . The corresponding density operator is:

$$(3.14) \quad \rho = |\psi\rangle\langle\psi| = \frac{1}{2}(|00\rangle\langle 00| + |00\rangle\langle 11| + |11\rangle\langle 00| + |11\rangle\langle 11|).$$

We apply the map  $\mathcal{I} \otimes \mathcal{T}$  (the partial transpose with respect to the second subsystem) to this state:

$$(3.15) \quad (\mathcal{I} \otimes \mathcal{T})[\rho] = \frac{1}{2}(|00\rangle\langle 00| + |01\rangle\langle 10| + |10\rangle\langle 01| + |11\rangle\langle 11|).$$

In the standard basis  $\{|00\rangle, |01\rangle, |10\rangle, |11\rangle\}$ , the matrix representation is

$$(3.16) \quad \frac{1}{2} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

This matrix has eigenvalues  $\{\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, -\frac{1}{2}\}$ . Since one eigenvalue is negative, the resulting operator is not positive. Thus,  $\mathcal{T}$  is not completely positive.  $\diamond$

We now arrive at the formal definition of a quantum channel.

**Definition 3.12** (Quantum channel, or CPTP map). *A **quantum channel**  $\mathcal{Q}$  is a linear map  $L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  that is completely positive (CP) and trace preserving (TP).*

If  $\mathcal{Q}$  is a quantum channel, it maps any density operator  $\rho \in \mathcal{D}(\mathbb{C}^N)$  to a density operator  $\mathcal{Q}[\rho] \in \mathcal{D}(\mathbb{C}^M)$ . The complete positivity condition ensures that if  $\mathcal{Q}$  acts locally on a subsystem of a larger entangled state  $\tilde{\rho} \in \mathcal{D}(\mathbb{C}^K \otimes \mathbb{C}^N)$ , the resulting state  $(\mathcal{I}_K \otimes \mathcal{Q})[\tilde{\rho}]$  remains a valid density operator in  $\mathcal{D}(\mathbb{C}^K \otimes \mathbb{C}^M)$ . This property is fundamental to the consistency of quantum mechanics.

**Example 3.13.** A fundamental class of quantum channels is the **unitary channel**. This requires the input and output dimensions to be equal,  $N = M$ . Given a unitary matrix  $U \in \mathbf{U}(N)$ , the corresponding channel  $\mathcal{U} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^N)$  acts by conjugation:

$$(3.17) \quad \mathcal{U}[\rho] = U\rho U^\dagger.$$

This map is trace-preserving, as  $\text{Tr}[U\rho U^\dagger] = \text{Tr}[\rho U^\dagger U] = \text{Tr}[\rho]$ . It is also completely positive, as we will see shortly. The identity channel  $\mathcal{I}$  is a unitary channel with  $U = I$ .  $\diamond$

A powerful way to characterize and construct quantum channels is through the Kraus representation.

**Proposition 3.14.** *Let  $\{K_j\}_{j \in [R]}$  be a set of matrices in  $\mathbb{C}^{M \times N}$  satisfying the completeness relation*

$$(3.18) \quad \sum_{j \in [R]} K_j^\dagger K_j = I_N.$$

*Then the linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  defined by*

$$(3.19) \quad \mathcal{Q}[\rho] = \sum_{j \in [R]} K_j \rho K_j^\dagger$$

*is a quantum channel (CPTP).*

**PROOF.** We first verify complete positivity. Let  $L$  be an arbitrary integer and consider any positive operator  $X \in L(\mathbb{C}^L \otimes \mathbb{C}^N)$ . The action of the extended map is

$$(3.20) \quad (\mathcal{I}_L \otimes \mathcal{Q})[X] = \sum_{j \in [R]} (I_L \otimes K_j) X (I_L \otimes K_j)^\dagger.$$

For any operator  $A$ , if  $X$  is positive, then  $AXA^\dagger$  is also positive. Thus, each term in the summation is a positive operator. Since the sum of positive operators is positive,  $\mathcal{I}_L \otimes \mathcal{Q}$  is a positive map for all  $L$ . Thus,  $\mathcal{Q}$  is completely positive.

Next, we verify the trace-preserving property. For any  $\rho \in L(\mathbb{C}^N)$ , using the linearity and the cyclic property of the trace, we have

$$(3.21) \quad \text{Tr}[\mathcal{Q}[\rho]] = \sum_{j \in [R]} \text{Tr}[K_j \rho K_j^\dagger] = \text{Tr} \left[ \rho \left( \sum_{j \in [R]} K_j^\dagger K_j \right) \right].$$

Substituting the completeness relation  $\sum_{j \in [R]} K_j^\dagger K_j = I_N$ , we obtain  $\text{Tr}[\rho I_N] = \text{Tr}[\rho]$ . Therefore,  $\mathcal{Q}$  is trace-preserving.  $\square$



The representation in Eq. (3.19) is called the **Kraus form** or the **operator sum representation** of the channel. The operators  $\{K_j\}$  are known as Kraus operators. For example, the unitary channel in Example 3.13 is in Kraus form with a single Kraus operator  $K_0 = U$ .

We can now explore the connection between classical stochastic evolution and quantum channels. This correspondence highlights that quantum mechanics is a generalization of classical probability theory.

For any probability distribution  $p \in \mathbb{R}^N$ , we can embed it into a quantum state

$$(3.22) \quad \rho = \sum_{i \in [N]} p_i |i\rangle\langle i|.$$

This diagonal density matrix is called a **classical state** or **probabilistic state**.

Given a (column) stochastic matrix  $P \in \mathbb{R}^{M \times N}$  (i.e.,  $P_{ij} \geq 0$  and  $\sum_{i \in [M]} P_{ij} = 1$  for all  $j \in [N]$ ), which defines a classical Markov process mapping distribution  $p$  to  $p' = Pp$ , we can construct a corresponding **classical channel**  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  defined by

$$(3.23) \quad \mathcal{Q}[\rho] = \sum_{i \in [M], j \in [N]} P_{ij} |i\rangle\langle j| \rho |j\rangle\langle i|.$$

If  $\rho$  is a classical state,  $\mathcal{Q}[\rho]$  is also a classical state corresponding to the evolved probability distribution  $p'$ .

**Exercise 3.1.** Prove that the classical channel  $\mathcal{Q}$  defined in Eq. (3.23) is indeed a quantum channel (CPTP).

The fact that classical channels are a subset of quantum channels suggests that any advantage offered by quantum computation must stem from the utilization of the off-diagonal entries of the density matrix (coherence) and the structure of non-classical channels.

We now present several examples of important quantum channels, typically modeling different types of noise processes in qubits ( $N = M = 2$ ).

**Example 3.15** (Bit flip and phase flip channels). The bit flip channel  $\mathcal{Q}_{\text{bf}}$  describes a process where the qubit state is flipped (i.e.,  $X$  gate applied) with probability  $1 - p$ , and remains unchanged with probability  $p$ :

$$(3.24) \quad \mathcal{Q}_{\text{bf}}[\rho] = p\rho + (1 - p)X\rho X, \quad 0 \leq p \leq 1.$$

This is in Kraus form with  $K_0 = \sqrt{p}I$  and  $K_1 = \sqrt{1 - p}X$ .

Similarly, the phase flip channel  $\mathcal{Q}_{\text{pf}}$  flips the relative phase (i.e.,  $Z$  gate applied) with probability  $1 - p$ :

$$(3.25) \quad \mathcal{Q}_{\text{pf}}[\rho] = p\rho + (1 - p)Z\rho Z, \quad 0 \leq p \leq 1.$$

This channel is also known as the **dephasing channel**, as it suppresses coherences while leaving populations unchanged.  $\diamond$

**Example 3.16** (Depolarizing channel). The depolarizing channel  $\mathcal{Q}_{\text{dp}} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^N)$  models a process where the state remains intact with probability  $p$ , and is replaced by the maximally mixed state  $I/N$  with probability  $1 - p$ :

$$(3.26) \quad \mathcal{Q}_{\text{dp}}[\rho] = p\rho + \frac{1 - p}{N}I, \quad 0 \leq p \leq 1.$$

$\diamond$

**Example 3.17** (Amplitude damping channel). The amplitude damping channel  $\mathcal{Q}_{\text{ad}} : L(\mathbb{C}^2) \rightarrow L(\mathbb{C}^2)$  models energy dissipation, such as spontaneous emission, where an excited state  $|1\rangle$  decays to the ground state  $|0\rangle$  with probability  $\gamma$ . It is described by the Kraus operators

$$(3.27) \quad K_0 = \begin{pmatrix} 1 & 0 \\ 0 & \sqrt{1-\gamma} \end{pmatrix}, \quad K_1 = \begin{pmatrix} 0 & \sqrt{\gamma} \\ 0 & 0 \end{pmatrix}, \quad 0 \leq \gamma \leq 1.$$

◇

Perhaps surprisingly, the converse of Proposition 3.14 is also true: every quantum channel can be written in the Kraus form. This fundamental result demonstrates that the abstract definition of a CPTP map is equivalent to the constructive definition provided by the operator sum representation.

**THEOREM 3.18** (Choi–Kraus Representation). *A linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  is a quantum channel if and only if there exists a set of matrices  $\{K_j\}_{j \in [R]}$  in  $\mathbb{C}^{M \times N}$ , with  $R \leq NM$ , satisfying the completeness relation  $\sum_{j \in [R]} K_j^\dagger K_j = I_N$ , such that  $\mathcal{Q}$  takes the form*

$$(3.28) \quad \mathcal{Q}(\rho) = \sum_{j \in [R]} K_j \rho K_j^\dagger.$$

**PROOF.** The “if” part is established by Proposition 3.14. We now prove the “only if” part using a technique known as the **Choi–Jamiołkowski isomorphism**.

Let  $\mathcal{Q}$  be a quantum channel. Define an unnormalized maximally entangled state on  $\mathbb{C}^N \otimes \mathbb{C}^N$ :

$$(3.29) \quad |\gamma\rangle = \sum_{i \in [N]} |i\rangle \otimes |i\rangle.$$

Let  $\mathcal{I}_N$  denote the identity map on the first  $N$ -dimensional register (the ancilla). By the complete positivity of  $\mathcal{Q}$ , the map  $\mathcal{I}_N \otimes \mathcal{Q}$  is positive. Therefore, the **Choi matrix** defined as

$$(3.30) \quad \sigma = (\mathcal{I}_N \otimes \mathcal{Q})[|\gamma\rangle\langle\gamma|] \in L(\mathbb{C}^N \otimes \mathbb{C}^M)$$

is a positive operator.

The Choi matrix completely characterizes the channel  $\mathcal{Q}$ . To see this, we use a key property of the maximally entangled state. For any vector  $|\psi\rangle = \sum_i \psi_i |i\rangle \in \mathbb{C}^N$ , let  $|\tilde{\psi}\rangle = \sum_i \bar{\psi}_i |i\rangle$  be its element-wise conjugate in the computational basis. We can verify the identity:

$$(3.31) \quad ((\langle\tilde{\psi}| \otimes I_N) |\gamma\rangle) = \sum_{i,j} \psi_j (\langle j|i\rangle \otimes |i\rangle) = \sum_i \psi_i |i\rangle = |\psi\rangle.$$

We can recover the action of the channel on  $|\psi\rangle\langle\psi|$  by taking the partial inner product of  $\sigma$  with  $|\tilde{\psi}\rangle$  on the first register. By the definition of the tensor product map and the identity above, we have:

$$(3.32) \quad \begin{aligned} ((\langle\tilde{\psi}| \otimes I_M) \sigma (|\tilde{\psi}\rangle \otimes I_M)) &= ((\langle\tilde{\psi}| \otimes I_M) (\mathcal{I}_N \otimes \mathcal{Q})[|\gamma\rangle\langle\gamma|] (|\tilde{\psi}\rangle \otimes I_M)) \\ &= \mathcal{Q} \left[ ((\langle\tilde{\psi}| \otimes I_N) |\gamma\rangle\langle\gamma| (|\tilde{\psi}\rangle \otimes I_N)) \right] \\ &= \mathcal{Q}(|\psi\rangle\langle\psi|). \end{aligned}$$

Since  $\sigma$  is positive, we can perform its eigendecomposition. Let  $R = \text{rank}(\sigma) \leq NM$ . We write

$$(3.33) \quad \sigma = \sum_{j \in [R]} |s_j\rangle\langle s_j|,$$

where  $|s_j\rangle \in \mathbb{C}^N \otimes \mathbb{C}^M$  are (potentially unnormalized) eigenvectors scaled by the square root of the eigenvalues.

For each  $j \in [R]$ , we define a linear operator  $K_j : \mathbb{C}^N \rightarrow \mathbb{C}^M$  via the relation (sometimes called vectorization or flattening):

$$(3.34) \quad K_j |\psi\rangle := (\langle \tilde{\psi} | \otimes I_M) |s_j\rangle.$$

Substituting the decomposition of  $\sigma$  back into the recovery formula:

$$(3.35) \quad \begin{aligned} \mathcal{Q}(|\psi\rangle\langle\psi|) &= (\langle \tilde{\psi} | \otimes I_M) \left( \sum_{j \in [R]} |s_j\rangle\langle s_j| \right) (|\tilde{\psi}\rangle \otimes I_M) \\ &= \sum_{j \in [R]} \left[ (\langle \tilde{\psi} | \otimes I_M) |s_j\rangle \right] \left[ \langle s_j | (|\tilde{\psi}\rangle \otimes I_M) \right] \\ &= \sum_{j \in [R]} (K_j |\psi\rangle)(K_j |\psi\rangle)^\dagger = \sum_{j \in [R]} K_j |\psi\rangle\langle\psi| K_j^\dagger. \end{aligned}$$

Since this holds for arbitrary  $|\psi\rangle$ , by linearity it holds for all operators  $\rho \in L(\mathbb{C}^N)$ .

Finally, we must verify the completeness relation. The trace-preserving property  $\text{Tr}[\mathcal{Q}(\rho)] = \text{Tr}[\rho]$  implies

$$(3.36) \quad \text{Tr} \left[ \sum_{j \in [R]} K_j \rho K_j^\dagger \right] = \text{Tr} \left[ \left( \sum_{j \in [R]} K_j^\dagger K_j \right) \rho \right] = \text{Tr}[I_N \rho].$$

Since this equality holds for all  $\rho$ , we must have  $\sum_{j \in [R]} K_j^\dagger K_j = I_N$ .  $\square$

The definition of complete positivity in Definition 3.10 requires verifying positivity for all dimensions  $K$ , which is operationally cumbersome. However, the proof of the Choi–Kraus theorem reveals that a much simpler criterion suffices. Let  $\mathcal{I}_N$  denote the identity channel on  $L(\mathbb{C}^N)$ . If we assume only that the map  $\mathcal{I}_N \otimes \mathcal{Q}$  is positive, then the Choi matrix  $\sigma$  (defined in the proof of Theorem 3.18) must be positive, as it is the image of the positive operator  $|\gamma\rangle\langle\gamma|$  under this map. As shown in the proof, the positivity of  $\sigma$  guarantees the existence of a Kraus representation for  $\mathcal{Q}$ . Finally, by Proposition 3.14, any map with a Kraus representation is completely positive (i.e.,  $\mathcal{I}_K \otimes \mathcal{Q}$  is positive for all  $K$ ). This establishes the equivalence between the original definition and a condition involving only an ancilla of the input dimension:

**Proposition 3.19** (Choi’s Theorem). *A linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  is completely positive if and only if its Choi matrix  $\sigma$  is positive semidefinite. Equivalently,  $\mathcal{Q}$  is CP if and only if the map  $\mathcal{I}_N \otimes \mathcal{Q}$  is positive.*

The Kraus representation provides deep insight into the structure of quantum channels. Another fundamental structural result is the Stinespring dilation theorem, which connects general quantum channels (which may involve decoherence or dissipation) to coherent evolution on a larger Hilbert space.

**THEOREM 3.20** (Stinespring dilation). *Given any quantum channel  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$ , there exists an ancilla system  $A$  of dimension  $R \leq NM$ , and an isometry  $V : \mathbb{C}^N \rightarrow \mathbb{C}^M \otimes \mathbb{C}^R$  (i.e.,  $V^\dagger V = I_N$ ) such that*

$$(3.37) \quad \mathcal{Q}(\rho) = \text{Tr}_A [V \rho V^\dagger].$$

Furthermore, this isometry can always be realized by a unitary evolution  $U$  on a sufficiently large joint system initialized with the ancilla in a fixed state  $|0\rangle$ :

$$(3.38) \quad \mathcal{Q}(\rho) = \text{Tr}_A [U(\rho \otimes |0\rangle\langle 0|)U^\dagger].$$

PROOF. By the Choi–Kraus theorem (Theorem 3.18),  $\mathcal{Q}$  has a Kraus representation  $\mathcal{Q}(\rho) = \sum_{j \in [R]} K_j \rho K_j^\dagger$ , where  $R \leq NM$ .

We construct the isometry  $V : \mathbb{C}^N \rightarrow \mathbb{C}^M \otimes \mathbb{C}^R$ . Let  $\{|j\rangle\}$  be an orthonormal basis for the ancilla space  $\mathbb{C}^R$ . Define  $V$  by

$$(3.39) \quad V|\psi\rangle = \sum_{j \in [R]} (K_j |\psi\rangle) \otimes |j\rangle.$$

We verify that  $V$  is an isometry. For any  $|\psi\rangle \in \mathbb{C}^N$ :

$$(3.40) \quad \begin{aligned} \langle \psi | V^\dagger V | \psi \rangle &= \|V|\psi\rangle\|^2 = \sum_{j \in [R]} \|K_j |\psi\rangle\|^2 \\ &= \sum_{j \in [R]} \langle \psi | K_j^\dagger K_j | \psi \rangle = \langle \psi | \left( \sum_j K_j^\dagger K_j \right) | \psi \rangle. \end{aligned}$$

By the completeness relation, this equals  $\langle \psi | \psi \rangle$ . Thus  $V^\dagger V = I_N$ .

Now we verify the representation in Eq. (3.37). We compute  $V\rho V^\dagger$ . It is helpful to view  $V$  formally as  $V = \sum_j K_j \otimes |j\rangle$ . Then

$$(3.41) \quad V\rho V^\dagger = \left( \sum_i K_i \otimes |i\rangle \right) \rho \left( \sum_j K_j^\dagger \otimes \langle j| \right) = \sum_{i,j} (K_i \rho K_j^\dagger) \otimes |i\rangle\langle j|.$$

Tracing over the ancilla (the second register) yields

$$(3.42) \quad \text{Tr}_A [V\rho V^\dagger] = \sum_{i,j} (K_i \rho K_j^\dagger) \text{Tr}[|i\rangle\langle j|] = \sum_j K_j \rho K_j^\dagger = \mathcal{Q}(\rho).$$

To realize this via a unitary evolution, we define  $U$  such that its action on the subspace corresponding to the initial state  $\rho \otimes |0\rangle\langle 0|$  matches the isometry  $V$ . Let the joint space be large enough (e.g., dimension  $D = \max(N, M)R$ ). We define  $U$  such that

$$(3.43) \quad U(|\psi\rangle \otimes |0\rangle) = V|\psi\rangle, \quad \forall |\psi\rangle \in \mathbb{C}^N.$$

(We might need to embed  $\mathbb{C}^N$  and  $\mathbb{C}^M \otimes \mathbb{C}^R$  into the larger space  $\mathbb{C}^D$ ). Since  $V$  is an isometry, this definition is norm-preserving. We can always extend this definition to a full unitary  $U$  on the joint space.

Finally, we verify the representation in Eq. (3.38). Let  $\rho = \sum_k p_k |\psi_k\rangle\langle \psi_k|$  be the spectral decomposition.

$$(3.44) \quad \begin{aligned} U(\rho \otimes |0\rangle\langle 0|)U^\dagger &= \sum_k p_k U(|\psi_k\rangle \otimes |0\rangle)(\langle \psi_k| \otimes \langle 0|)U^\dagger \\ &= \sum_k p_k (V|\psi_k\rangle)(V|\psi_k\rangle)^\dagger = V\rho V^\dagger. \end{aligned}$$

Therefore,  $\mathcal{Q}(\rho) = \text{Tr}_A [V\rho V^\dagger] = \text{Tr}_A [U(\rho \otimes |0\rangle\langle 0|)U^\dagger]$ . □

Theorem 3.20 states that any quantum channel, no matter how noisy or irreversible it appears, can always be modeled as a unitary interaction between the system and an environment (ancilla), followed by discarding the environment. This provides a powerful conceptual tool, showing that all quantum evolution is fundamentally unitary if we consider a large enough closed system.

### 3.3. Distance between state vectors and unitaries

A **distance** (also called a **metric**) on a set  $X$  is a function  $d : X \times X \rightarrow \mathbb{R}$  that assigns a real number  $d(x, y)$  to each pair of points  $x, y \in X$ . This function satisfies the following properties for all  $x, y, z \in X$ :

- (1) (Non-negativity)  $d(x, y) \geq 0$ .
- (2) (Identity of indiscernibles)  $d(x, y) = 0$  if and only if  $x = y$ .
- (3) (Symmetry)  $d(x, y) = d(y, x)$ .
- (4) (Triangle inequality)  $d(x, y) \leq d(x, z) + d(z, y)$ .

For example, the vector 2-norm defines a metric on  $\mathbb{C}^N : (x, y) \rightarrow \|x - y\|$ , and the operator norm defines a metric on  $U(N) : (U, V) \rightarrow \|U - V\|$ .

The difference for the product of  $K$  unitaries can be bounded using a simple technique sometimes referred to as a “hybrid argument”. This technique is used to bound the distance between two states by considering a sequence of “hybrid” unitaries, each of which differs from the next in the sequence by a small amount.

**Proposition 3.21** (Linear error growth for products of unitaries). *Given unitaries  $U_1, \tilde{U}_1, \dots, U_K, \tilde{U}_K \in U(N)$  satisfying*

$$(3.45) \quad \|U_i - \tilde{U}_i\| \leq \epsilon, \quad \forall i = 1, \dots, K,$$

*we have*

$$(3.46) \quad \|U_K \cdots U_1 - \tilde{U}_K \cdots \tilde{U}_1\| \leq K\epsilon.$$

PROOF. Use a telescoping series

$$(3.47) \quad \begin{aligned} & U_K \cdots U_1 - \tilde{U}_K \cdots \tilde{U}_1 \\ &= (U_K \cdots U_2 U_1 - U_K \cdots U_2 \tilde{U}_1) + (U_K \cdots U_3 U_2 \tilde{U}_1 - U_K \cdots U_3 \tilde{U}_2 \tilde{U}_1) + \cdots \\ & \quad + (U_K U_{K-1} \cdots \tilde{U}_1 - \tilde{U}_K \tilde{U}_{K-1} \cdots \tilde{U}_1) \\ &= U_K \cdots U_2 (U_1 - \tilde{U}_1) + U_K \cdots U_3 (U_2 - \tilde{U}_2) \tilde{U}_1 + \cdots + (U_K - \tilde{U}_K) \tilde{U}_{K-1} \cdots \tilde{U}_1. \end{aligned}$$

Since all  $U_i, \tilde{U}_i$  are unitary matrices, we readily have

$$(3.48) \quad \|U_K \cdots U_1 - \tilde{U}_K \cdots \tilde{U}_1\| \leq \sum_{i=1}^K \|U_i - \tilde{U}_i\| \leq K\epsilon.$$

□

For most of this book, the vector 2-norm and the operator norm distances are both convenient and sufficient. However, they are only applicable to pure states. For measuring the distance between mixed states, new tools will be needed. Even for pure states, unitaries may differ by a phase which should be inconsequential for measuring physical observables. These require the introduction of new metrics.

Two state vectors  $|\psi\rangle, |\varphi\rangle \in \mathbb{C}^N$  are physically indistinguishable if they only differ by a global phase. Similarly, two unitary matrices  $U, V \in \text{U}(N)$  induce the same evolution on density operators if they only differ by a global phase. Consider the matrices

$$(3.49) \quad I_+ := \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad I_- := \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix}.$$

We have in this case that  $\|I_+ - I_-\| = 2$ . However, for an arbitrary density matrix  $\rho$ , the induced evolution of the density operator under these two operators is

$$(3.50) \quad \|I_+ \rho I_+ - I_- \rho I_-\| = \|\rho - (-1)^2 \rho\| = 0.$$

This motivates the definition of the global phase invariant distance for vectors and unitary matrices. The subscript  $p$  in  $D_p$  stands for phase.

**Definition 3.22.** Let  $|\psi\rangle, |\varphi\rangle \in \mathbb{C}^N$  be two state vectors, their **global phase invariant distance** is

$$(3.51) \quad D_p(|\psi\rangle, |\varphi\rangle) := \min_{\phi \in \mathbb{R}} \|\psi\rangle - e^{i\phi} |\varphi\rangle\|.$$

**Definition 3.23.** For two unitaries  $U, V \in \text{U}(N)$ , their **global phase invariant distance** is

$$(3.52) \quad D_p(U, V) = \min_{\phi \in \mathbb{R}} \|U - e^{i\phi} V\|.$$

An **equivalence relation** on a set  $X$  is a binary relation  $\sim$  that satisfies the following three properties for all  $a, b, c \in X$ :

- (1) (Reflexivity)  $a \sim a$ .
- (2) (Symmetry) If  $a \sim b$ , then  $b \sim a$ .
- (3) (Transitivity) If  $a \sim b$  and  $b \sim c$ , then  $a \sim c$ .

A relation that satisfies these properties is called an equivalence relation, and it partitions the set  $X$  into disjoint equivalence classes.

**Definition 3.24.** Let  $X$  be a set and  $\sim$  be an equivalence relation on  $X$ . The quotient space (or quotient set)  $X/\sim$  is defined as the set of equivalence classes of  $X$  under the relation  $\sim$ . An equivalence class  $[x]$  of an element  $x \in X$  is the set of all elements in  $X$  that are equivalent to  $x$ , i.e.,

$$(3.53) \quad [x] = \{y \in X \mid y \sim x\}.$$

The quotient space  $X/\sim$  is the set of all such equivalence classes:

$$(3.54) \quad X/\sim = \{[x] \mid x \in X\}.$$

**Example 3.25.** Define an equivalence relation on  $\mathbb{C}^N$ :

$$(3.55) \quad x \sim y \iff x = \lambda y \text{ for some } \lambda \in \mathbb{C} \setminus \{0\}, \quad x, y \in \mathbb{C}^N \setminus \{0\}.$$

Then  $\text{PC}^N := \mathbb{C}^N \setminus \{0\} / \sim$  is called the **complex projective space**, which is isomorphic to the set of all nonzero physical states. The **real dimension** of a manifold  $M$  is the number of real coordinates needed to locally describe the manifold. For example, the real dimension of  $\mathbb{C}^N$  is  $2N$ , and the real dimension of  $\text{PC}^N$  is  $2N - 2$ .

We may identify each single qubit quantum state with a unique point on the Bloch sphere as

$$(3.56) \quad \mathbf{a} = (\sin \theta \cos \varphi, \sin \theta \sin \varphi, \cos \theta)^\top, \quad \theta, \varphi \in \mathbb{R}.$$

This agrees with the previous statement that the real dimension of  $\mathbb{P}\mathbb{C}^2$  is 2.  $\diamond$

**Exercise 3.2.** Prove that the global phase invariant distance is a distance on the complex projective space  $\mathbb{P}\mathbb{C}^N$ .

**Example 3.26.** Define an equivalence relation on  $U(N)$ :

$$(3.57) \quad U \sim V \iff U = e^{i\theta}V \text{ for some } \theta \in \mathbb{R}, \quad U, V \in U(N).$$

Then  $\text{PU}(N) := U(N)/\sim$  is called the **projective unitary group**. The real dimension of  $U(N)$  is  $N^2$ , and the real dimension of  $\text{PU}(N)$  is  $N^2 - 1$ .

Recall that the special unitary group  $\text{SU}(N)$  consists of all unitary matrices with determinant 1. So the real dimension of  $\text{SU}(N)$  is  $N^2 - 1$ . However, the equivalence relation on  $\text{SU}(N)$  is

$$(3.58) \quad U \sim V \iff U = e^{i2\pi k/N}V \text{ for some } k \in [N], \quad U, V \in \text{SU}(N).$$

So each equivalence class only consists of  $N$  discrete elements and does not reduce the dimension. Therefore the real dimension of the projective special unitary group denoted by  $\text{PSU}(N)$  is still  $N^2 - 1$ .  $\diamond$

**Exercise 3.3.** Prove that the global phase invariant distance is a distance on the projective unitary group  $\text{PU}(N)$ .

**Exercise 3.4.** Given unitaries  $U_1, \tilde{U}_1, \dots, U_K, \tilde{U}_K \in U(N)$  satisfying

$$(3.59) \quad D_p(U_i, \tilde{U}_i) \leq \epsilon, \quad \forall i = 1, \dots, K,$$

prove that

$$(3.60) \quad D_p(U_K \cdots U_1, \tilde{U}_K \cdots \tilde{U}_1) \leq K\epsilon.$$

Let  $|\varphi\rangle = e^{i\alpha} \cos \theta |\psi\rangle + \sin \theta |\perp\rangle$ , where  $\langle \psi | \perp \rangle = 0$  and  $0 \leq \theta \leq \pi/2$ . Then  $\cos \theta = |\langle \varphi | \psi \rangle|$  is the overlap between the two vectors. We can perform a unitary operation that rotates  $e^{i\alpha} |\psi\rangle$  to  $|0\rangle$  and  $|\perp\rangle$  to  $|1\rangle$ . Direct calculation shows

$$(3.61) \quad D_p(|\psi\rangle, |\varphi\rangle) = \min_{\phi \in \mathbb{R}} \|\lvert 0 \rangle - e^{i\phi}(e^{i\alpha} \cos \theta \lvert 0 \rangle + \sin \theta \lvert 1 \rangle)\| = \sqrt{2(1 - \cos \theta)} = \sqrt{2(1 - |\langle \varphi | \psi \rangle|)}.$$

Therefore the global phase invariant distance between two vectors can be directly computed from the overlap.

**Exercise 3.5.** For  $U, V \in U(N)$ , prove that

$$(3.62) \quad D_p(U, V) = 2 \min_{\phi} \max_j \left| \sin \frac{\lambda_j - \phi}{2} \right|,$$

where  $\{e^{i\lambda_j}\}$  are eigenvalues of  $V^\dagger U$ .

**Exercise 3.6.** For  $U, V \in U(N)$ , another distance that is invariant to the global phase is

$$(3.63) \quad D_{p,F}(U, V) = \frac{1}{\sqrt{2N}} \min_{\phi} \|U - e^{i\phi}V\|_F.$$

Prove that

$$(3.64) \quad D_{p,F}(U, V) = \sqrt{1 - \frac{|\text{Tr}[U^\dagger V]|}{N}}.$$

### 3.4. Distance between classical states and classical channels

In this section, we provide a connection between concepts in classical probabilistic computation and density operators and quantum channels in quantum computation. For two probability distributions  $p, q \in \mathbb{R}^N$ , the **total variation distance** is

$$(3.65) \quad D(p, q) := \frac{1}{2} \sum_{i \in [N]} |p_i - q_i|.$$

The name total variation distance comes from that it measures the largest difference between  $p$  and  $q$  for some subset (also called event)  $S$ . The total variation distance is the default metric we will use between probability distributions and will be denoted by  $D$  without subscripts.

**Proposition 3.27.** *For any two classical probability distributions  $p, q \in \mathbb{R}^N$ ,*

$$(3.66) \quad D(p, q) = \max_S (p(S) - q(S)) := \max_S \left( \sum_{i \in S} p_i - \sum_{i \in S} q_i \right),$$

where the maximization is over all subsets  $S$ .

PROOF. For any subset  $S$ , let  $\bar{S}$  be its complement. Then

$$(3.67) \quad 0 = \sum_i p_i - \sum_i q_i = \sum_{i \in S} (p_i - q_i) + \sum_{i \in \bar{S}} (p_i - q_i).$$

Hence

$$(3.68) \quad \sum_{i \in S} p_i - \sum_{i \in S} q_i = \frac{1}{2} \left( \sum_{i \in S} (p_i - q_i) - \sum_{i \in \bar{S}} (p_i - q_i) \right) \leq D(p, q).$$

Now let  $S = \{i | p_i \geq q_i\}$ . Then

$$(3.69) \quad \frac{1}{2} \left( \sum_{i \in S} (p_i - q_i) - \sum_{i \in \bar{S}} (p_i - q_i) \right) = \frac{1}{2} \sum_i |p_i - q_i| = D(p, q),$$

and the equality is achieved.  $\square$

We now prove that the application of a transition matrix does not increase the total variation distance.

**Proposition 3.28.** *Given a transition matrix  $P \in \mathbb{R}^{N \times N}$ , and any two classical probability distributions  $p, q \in \mathbb{R}^N$ ,*

$$(3.70) \quad D(Pp, Pq) \leq D(p, q).$$

*If the equality holds for any  $p, q \in \mathbb{R}^N$ , then  $P$  is a permutation matrix.*

PROOF. Use the left stochasticity of the transition matrix, we have

$$(3.71) \quad D(Pp, Pq) = \frac{1}{2} \sum_i \left| \sum_j P_{ij} (p_j - q_j) \right| \leq \frac{1}{2} \sum_i \sum_j P_{ij} |p_j - q_j| = \frac{1}{2} \sum_j |p_j - q_j| = D(p, q).$$



If the equality holds for any  $p, q \in \mathbb{R}^N$ , we prove that each row of  $P$  has only one nonzero entry. If this is not the case, assume that there exists a row index  $i$  and two distinct column indices  $j_1 \neq j_2$  such that  $P_{ij_1} > 0$  and  $P_{ij_2} > 0$ . Choose  $p = e_{j_1}$  and  $q = e_{j_2}$ . Then for this row  $i$ ,

$$(3.72) \quad \left| \sum_j P_{ij}(p_j - q_j) \right| = |P_{i,j_1} - P_{i,j_2}| < P_{i,j_1} + P_{i,j_2} = \sum_j P_{ij} |p_j - q_j|,$$

which contradicts equality in the triangle inequality step above. Hence each row has exactly one nonzero entry. By left stochasticity, each column must also have exactly one nonzero entry, which must equal 1. This proves that  $P$  is a permutation matrix.  $\square$

The **induced total variation distance** between two transition matrices  $P, Q \in \mathbb{R}^{N \times N}$  is defined as

$$(3.73) \quad D(P, Q) = \max_{j \in [N]} D(P_{:,j}, Q_{:,j}).$$

**Exercise 3.7.** Prove that  $D(\cdot, \cdot)$  is a distance on the set of  $N \times N$  transition matrices.

Finally, we prove that the difference for the composition of  $K$  classical channels grows linearly.

**Proposition 3.29** (Linear error growth for product of transition matrices). *Given the transition matrices  $P_1, \tilde{P}_1, \dots, P_K, \tilde{P}_K \in \mathbb{R}^{N \times N}$ , the induced total variation distance satisfies*

$$(3.74) \quad D(P_K \cdots P_1, \tilde{P}_K \cdots \tilde{P}_1) \leq \sum_{i=1}^K D(P_i, \tilde{P}_i).$$

PROOF. Using the telescope series Proposition 3.21, it is sufficient to consider the case for  $K = 2$ . Then

$$(3.75) \quad \begin{aligned} D(P_2 P_1, \tilde{P}_2 \tilde{P}_1) &\leq D(P_2 P_1, P_2 \tilde{P}_1) + D(P_2 \tilde{P}_1, \tilde{P}_2 \tilde{P}_1) \\ &= \max_{j \in [N]} D((P_2 P_1)_{:,j}, (P_2 \tilde{P}_1)_{:,j}) + \max_{j \in [N]} D((P_2 \tilde{P}_1)_{:,j}, (\tilde{P}_2 \tilde{P}_1)_{:,j}) \\ &\leq \max_{j \in [N]} D((P_1)_{:,j}, (\tilde{P}_1)_{:,j}) + \max_{j \in [N]} \left( \max_{l \in [N]} D((P_2)_{:,l}, (\tilde{P}_2)_{:,l}) \right) \sum_k (\tilde{P}_1)_{kj} \\ &\leq \max_{j \in [N]} D((P_1)_{:,j}, (\tilde{P}_1)_{:,j}) + \max_{l \in [N]} D((P_2)_{:,l}, (\tilde{P}_2)_{:,l}) \\ &= D(P_1, \tilde{P}_1) + D(P_2, \tilde{P}_2). \end{aligned}$$

Here we have used Proposition 3.28 and the left stochasticity of  $\tilde{P}_1$ .  $\square$

### 3.5. Distance between quantum states

Quantifying the similarity or difference between quantum states is fundamental to quantum information theory. It allows us to analyze the performance of quantum algorithms, assess the errors in quantum communication protocols, and understand the distinguishability of quantum states through measurements. In this section, we introduce the two most widely used measures: the trace distance and the fidelity. These generalize the corresponding concepts for classical probability distributions, such as the total variation distance discussed in Section 3.4. For a comprehensive treatment, we refer readers to [NC00, Chapter 9] and [Wat18, Chapter 3].

**3.5.1. Schatten norms and the trace norm.** To define distances between density operators, which are matrices, we first need appropriate matrix norms. The Schatten norms provide a family of norms generalizing the  $\ell^p$  norms for vectors to the space of operators.

Let  $A \in \mathbb{C}^{M \times N}$ . The singular values of  $A$ , denoted  $\sigma_i(A)$ , are the square roots of the non-negative eigenvalues of  $A^\dagger A$ . The Schatten  $p$ -norm of  $A$  for  $p \geq 1$  is defined as the  $\ell^p$  norm of its singular values:

$$(3.76) \quad \|A\|_p := \left( \sum_i \sigma_i(A)^p \right)^{\frac{1}{p}}.$$

This can also be expressed using the trace function. Let  $|A| := \sqrt{A^\dagger A}$  denote the positive semidefinite square root of  $A^\dagger A$ . Then

$$(3.77) \quad \|A\|_p = (\text{Tr}[|A|^p])^{\frac{1}{p}}.$$

The following choices of  $p$  are particularly important:

- The Schatten 1-norm, also known as the trace norm, is the sum of the singular values:

$$(3.78) \quad \|A\|_1 = \text{Tr}[|A|] = \sum_i \sigma_i(A).$$

If  $A$  is positive semidefinite,  $|A| = A$ , so  $\|A\|_1 = \text{Tr}[A]$ .

- The Schatten 2-norm (also called the Hilbert-Schmidt norm or Frobenius norm) is the Euclidean norm of the singular values:

$$(3.79) \quad \|A\|_2 = \sqrt{\text{Tr}[A^\dagger A]} = \left( \sum_i \sigma_i(A)^2 \right)^{\frac{1}{2}}.$$

- The Schatten  $\infty$ -norm is the maximum singular value:

$$(3.80) \quad \|A\|_\infty = \lim_{p \rightarrow \infty} \|A\|_p = \max_i \sigma_i(A).$$

This is identical to the standard operator norm (the induced  $\ell^2 \rightarrow \ell^2$  norm), often denoted  $\|A\|$  (equivalently  $\|A\|_\infty$ ).

A basic but useful property relates the trace of a matrix to its trace norm.

**Proposition 3.30.** *For any square matrix  $A \in L(\mathbb{C}^N)$ ,*

$$(3.81) \quad |\text{Tr}[A]| \leq \|A\|_1.$$

PROOF. Consider the singular value decomposition  $A = U\Sigma V^\dagger$ , where  $U, V$  are unitary and  $\Sigma = \text{diag}(\sigma_i)$  contains the singular values. Using the cyclic property of the trace:

$$(3.82) \quad \text{Tr}[A] = \text{Tr}[U\Sigma V^\dagger] = \text{Tr}[\Sigma V^\dagger U].$$

Let  $W = V^\dagger U$ . Since  $W$  is unitary, its entries satisfy  $|W_{ii}| \leq 1$  for all  $i$ . Therefore, by the triangle inequality,

$$(3.83) \quad |\text{Tr}[A]| = \left| \sum_i \sigma_i W_{ii} \right| \leq \sum_i \sigma_i |W_{ii}| \leq \sum_i \sigma_i = \|A\|_1.$$

□

The Schatten norms share many properties with the  $\ell^p$  norms for vectors, including the triangle inequality and Hölder's inequality. We state these fundamental results without proof, referring the reader to texts on matrix analysis such as [Bha97].

**Proposition 3.31** (Properties of Schatten  $p$ -norms). *Let  $A, B$  be operators.*

- (1) (*Triangle inequality*) For  $1 \leq p \leq \infty$ ,  $\|A + B\|_p \leq \|A\|_p + \|B\|_p$ .
- (2) (*Hölder's inequality*, [Bha97, Corollary IV.2.6]) For  $1 \leq p, q \leq \infty$  satisfying  $\frac{1}{p} + \frac{1}{q} = 1$ , if the product  $AB$  is defined, then  $\|AB\|_1 \leq \|A\|_p \|B\|_q$ .

We are primarily interested in the trace norm ( $p = 1$ ) and the operator norm ( $p = \infty$ ). An important specialization of Hölder's inequality is the case  $p = \infty, q = 1$ :

$$(3.84) \quad \|AB\|_1 \leq \|A\|_\infty \|B\|_1.$$

This inequality is frequently used to bound the trace norm of a product. Another useful variation involves the trace of a product, which can be viewed as a generalization of the Cauchy-Schwarz inequality. We provide a self-contained proof of this specific case.

**Lemma 3.32** (Hölder's inequality for trace). *For any operators  $A, B \in L(\mathbb{C}^N)$ , the following inequality holds:*

$$(3.85) \quad |\mathrm{Tr}(A^\dagger B)| \leq \|A\|_\infty \|B\|_1.$$

PROOF. Let  $B = U\Sigma V^\dagger$  be the SVD of  $B$ , with singular values  $s_i$ . By definition,  $\|B\|_1 = \sum_i s_i$ . Using the cyclic property of the trace:

$$(3.86) \quad \mathrm{Tr}(A^\dagger B) = \mathrm{Tr}(A^\dagger U \Sigma V^\dagger) = \mathrm{Tr}(V^\dagger A^\dagger U \Sigma).$$

Let  $W = V^\dagger A^\dagger U$ . Since  $U$  and  $V$  are unitary, the operator norm is invariant under unitary multiplication:  $\|W\|_\infty = \|A^\dagger\|_\infty$ . Furthermore,  $\|A^\dagger\|_\infty = \|A\|_\infty$  as they share the same singular values. The trace is the sum of the diagonal elements of  $W$  weighted by the singular values:

$$(3.87) \quad \mathrm{Tr}(W \Sigma) = \sum_i W_{ii} s_i.$$

We can now bound the magnitude of the trace using the triangle inequality:

$$(3.88) \quad \begin{aligned} |\mathrm{Tr}(A^\dagger B)| &= \left| \sum_i W_{ii} s_i \right| \leq \sum_i |W_{ii}| s_i \\ &\leq \sum_i \|W\|_\infty s_i = \|A\|_\infty \sum_i s_i = \|A\|_\infty \|B\|_1. \end{aligned}$$

□

We now consider how the trace norm behaves under the partial trace operation, which often arises when dealing with composite systems.

**Exercise 3.8.** Let  $|u\rangle, |v\rangle$  be normalized state vectors in  $\mathcal{H}_A \otimes \mathcal{H}_B$ . Show that

$$(3.89) \quad \|\mathrm{Tr}_B |u\rangle\langle v|\|_1 \leq 1.$$

(Hint: use Hölder's inequality for the Schatten 2-norm.)

More generally, the partial trace is a contraction with respect to the trace norm.

**Proposition 3.33** (Partial trace does not increase the trace norm). *For any operator  $O \in L(\mathcal{H}_A \otimes \mathcal{H}_B)$ ,*

$$(3.90) \quad \|\mathrm{Tr}_B O\|_1 \leq \|O\|_1.$$

PROOF. Consider the singular value decomposition of the operator  $O$ :

$$(3.91) \quad O = \sum_k \sigma_k |u_k\rangle\langle v_k|,$$

where  $\sigma_k > 0$  are the singular values, and  $\{|u_k\rangle\}, \{|v_k\rangle\}$  are sets of orthonormal vectors in  $\mathcal{H}_A \otimes \mathcal{H}_B$ . The trace norm is  $\|O\|_1 = \sum_k \sigma_k$ .

Applying the partial trace and using the triangle inequality (Proposition 3.31):

$$(3.92) \quad \|\mathrm{Tr}_B O\|_1 = \left\| \sum_k \sigma_k \mathrm{Tr}_B |u_k\rangle\langle v_k| \right\|_1 \leq \sum_k \sigma_k \|\mathrm{Tr}_B |u_k\rangle\langle v_k|\|_1.$$

By Exercise 3.8,  $\|\mathrm{Tr}_B |u_k\rangle\langle v_k|\|_1 \leq 1$ . Therefore,

$$(3.93) \quad \|\mathrm{Tr}_B O\|_1 \leq \sum_k \sigma_k = \|O\|_1.$$

□

The trace norm and the operator norm are dual to each other with respect to the trace inner product, a property that is frequently exploited in optimization problems and for deriving operational interpretations of these norms.

**Lemma 3.34** (Duality of Trace and Operator Norms). *For any operator  $Y \in L(\mathbb{C}^N)$ , the following identities hold:*

$$(3.94) \quad \|Y\|_1 = \sup_{\|Z\|_\infty \leq 1} |\mathrm{Tr}(Z^\dagger Y)|,$$

and

$$(3.95) \quad \|Y\|_\infty = \sup_{\|X\|_1 \leq 1} |\mathrm{Tr}(Y^\dagger X)|.$$

PROOF. We first prove Eq. (3.94). Let  $S_1$  denote the right-hand side. Applying Hölder's inequality (Lemma 3.32), we have  $|\mathrm{Tr}(Z^\dagger Y)| \leq \|Z\|_\infty \|Y\|_1$ . If we restrict the optimization to  $\|Z\|_\infty \leq 1$ , then  $|\mathrm{Tr}(Z^\dagger Y)| \leq \|Y\|_1$ . Taking the supremum yields  $S_1 \leq \|Y\|_1$ .

To show  $S_1 \geq \|Y\|_1$ , we construct an operator  $Z$  that achieves the bound. Let  $Y = U\Sigma V^\dagger$  be the SVD of  $Y$ . Define  $Z = UV^\dagger$ . Since  $Z$  is unitary,  $\|Z\|_\infty = 1$ . We compute the trace:

$$(3.96) \quad \begin{aligned} \mathrm{Tr}(Z^\dagger Y) &= \mathrm{Tr}((UV^\dagger)^\dagger (U\Sigma V^\dagger)) = \mathrm{Tr}(VU^\dagger U\Sigma V^\dagger) \\ &= \mathrm{Tr}(V\Sigma V^\dagger) = \mathrm{Tr}(\Sigma) = \|Y\|_1. \end{aligned}$$

Thus,  $S_1 \geq \|Y\|_1$ .

Next, we prove Eq. (3.95). Let  $S_\infty$  denote the right-hand side. Applying Lemma 3.32, we have  $|\mathrm{Tr}(Y^\dagger X)| \leq \|Y\|_\infty \|X\|_1$ . Restricting to  $\|X\|_1 \leq 1$  and taking the supremum yields  $S_\infty \leq \|Y\|_\infty$ .

To show  $S_\infty \geq \|Y\|_\infty$ , we construct an optimal  $X$ . Let  $Y = \sum_i s_i |u_i\rangle\langle v_i|$  be the SVD of  $Y$ , ordered such that  $s_1 = \|Y\|_\infty$ . Define the rank-1 operator  $X = |u_1\rangle\langle v_1|$ . Since  $|u_1\rangle, |v_1\rangle$  are

normalized,  $\|X\|_1 = 1$ . We compute the trace:

$$\begin{aligned} \text{Tr}(Y^\dagger X) &= \text{Tr} \left( \left( \sum_i s_i |v_i\rangle\langle u_i| \right) |u_1\rangle\langle v_1| \right) \\ (3.97) \quad &= \text{Tr} \left( \sum_i s_i |v_i\rangle\langle v_1| \langle u_i|u_1\rangle \right). \end{aligned}$$

Due to the orthonormality of  $\{|u_i\rangle\}$ , only the  $i = 1$  term survives:

$$(3.98) \quad \text{Tr}(Y^\dagger X) = \text{Tr}(s_1 |v_1\rangle\langle v_1|) = s_1 = \|Y\|_\infty.$$

Thus,  $S_\infty \geq \|Y\|_\infty$ .  $\square$

When the operator  $Y$  is Hermitian, the optimization domains in these duality relations can also be restricted to Hermitian operators.

**Lemma 3.35** (Duality for Hermitian Operators). *Let  $H \in L(\mathbb{C}^N)$  be a Hermitian operator.*

- (1) *The trace norm is achieved by maximizing over Hermitian operators in the unit operator-norm ball (i.e.,  $-I \preceq Z \preceq I$ ):*

$$(3.99) \quad \|H\|_1 = \sup\{|\text{Tr}(ZH)| : Z = Z^\dagger, \|Z\|_\infty \leq 1\}.$$

- (2) *The operator norm is achieved by maximizing over density operators:*

$$(3.100) \quad \|H\|_\infty = \sup\{|\text{Tr}(H\rho)| : \rho \in \mathcal{D}(\mathbb{C}^N)\}.$$

PROOF. In both cases, the inequality  $\leq$  (for the left-hand side) follows immediately from Lemma 3.34, as the restricted optimization domains are subsets of the original domains. We only need to show that the bounds can be achieved within these restricted domains.

1. Proof of Eq. (3.99). Let  $H = \sum_i \lambda_i |\psi_i\rangle\langle\psi_i|$  be the spectral decomposition, where  $\lambda_i \in \mathbb{R}$ . The trace norm is  $\|H\|_1 = \sum_i |\lambda_i|$ . Define the sign operator  $Z = \sum_i \text{sgn}(\lambda_i) |\psi_i\rangle\langle\psi_i|$ .  $Z$  is Hermitian, and its eigenvalues are in  $\{-1, 0, 1\}$ , so  $\|Z\|_\infty \leq 1$ .

$$(3.101) \quad \text{Tr}(ZH) = \sum_i \text{sgn}(\lambda_i) \lambda_i = \sum_i |\lambda_i| = \|H\|_1.$$

2. Proof of Eq. (3.100). The operator norm is  $\|H\|_\infty = \max_i |\lambda_i|$ . Let  $k$  be an index achieving the maximum. Define the pure state  $\rho = |\psi_k\rangle\langle\psi_k|$ , which is a density operator.

$$(3.102) \quad |\text{Tr}(H\rho)| = |\langle\psi_k|H|\psi_k\rangle| = |\lambda_k| = \|H\|_\infty. \quad \square$$

**3.5.2. Trace distance.** The trace norm provides a natural way to define a distance metric on the space of quantum states, generalizing the classical total variation distance.

**Definition 3.36** (Trace distance). *The **trace distance** between two quantum states  $\rho, \sigma \in \mathcal{D}(\mathbb{C}^N)$  is defined as*

$$(3.103) \quad D(\rho, \sigma) := \frac{1}{2} \|\rho - \sigma\|_1.$$

The factor of  $1/2$  ensures that the distance lies in the range  $[0, 1]$ . Since  $\|\rho\|_1 = 1$  and  $\|\sigma\|_1 = 1$ , the triangle inequality (Proposition 3.31) gives  $\|\rho - \sigma\|_1 \leq \|\rho\|_1 + \|\sigma\|_1 = 2$ .

**Example 3.37** (Trace distance for classical states). Consider classical probability distributions  $p, s \in \mathbb{R}^N$  embedded as classical states:

$$(3.104) \quad \rho = \sum_{i \in [N]} p_i |i\rangle\langle i|, \quad \sigma = \sum_{i \in [N]} s_i |i\rangle\langle i|.$$

The difference  $\rho - \sigma$  is a diagonal matrix with entries  $p_i - s_i$ . The trace norm is the sum of the absolute values of the eigenvalues:

$$(3.105) \quad D(\rho, \sigma) = \frac{1}{2} \|\rho - \sigma\|_1 = \frac{1}{2} \sum_i |p_i - s_i|.$$

This is exactly the total variation distance  $D(p, s)$  between the probability distributions  $p$  and  $s$ .  $\diamond$

The trace distance has an operational interpretation related to the distinguishability of quantum states through measurement. This is the quantum generalization of Proposition 3.27.

**Proposition 3.38** (Operational interpretation of trace distance). *For any quantum states  $\rho, \sigma \in \mathcal{D}(\mathbb{C}^N)$ , the trace distance satisfies*

$$(3.106) \quad D(\rho, \sigma) = \max_{0 \preceq M \preceq I} \text{Tr}[M(\rho - \sigma)].$$

*The maximum is achieved when  $M$  is the projector onto the subspace where  $\rho - \sigma$  is positive.*

PROOF. Let  $\Delta = \rho - \sigma$ .  $\Delta$  is Hermitian and  $\text{Tr}[\Delta] = 0$ . We want to maximize  $\text{Tr}[M\Delta]$  over  $0 \preceq M \preceq I$ .

We utilize the duality results established earlier. Consider an operator  $M$  such that  $0 \preceq M \preceq I$ . Define  $Z = 2M - I$ . Then  $Z$  is Hermitian, and  $-I \preceq Z \preceq I$ , which means  $\|Z\|_\infty \leq 1$ . We have

$$(3.107) \quad \text{Tr}[Z\Delta] = \text{Tr}[(2M - I)\Delta] = 2\text{Tr}[M\Delta].$$

By the Hermitian duality relation (Lemma 3.35, Eq. (3.99)),  $\|\Delta\|_1 = \sup\{|\text{Tr}(Z'\Delta)| : Z' = Z'^\dagger, \|Z'\|_\infty \leq 1\}$ . Since  $Z$  is admissible for this optimization, we have

$$(3.108) \quad 2\text{Tr}[M\Delta] = \text{Tr}[Z\Delta] \leq \|\Delta\|_1.$$

Thus,  $\text{Tr}[M\Delta] \leq \frac{1}{2} \|\Delta\|_1 = D(\rho, \sigma)$ .

To show equality, we construct an optimal  $M$ . Let  $\Delta = \Delta_+ - \Delta_-$ , where  $\Delta_+, \Delta_-$  are positive semidefinite operators with orthogonal support. Since  $\text{Tr}[\Delta] = 0$ , we have  $\text{Tr}[\Delta_+] = \text{Tr}[\Delta_-]$ . The trace norm is

$$(3.109) \quad \|\Delta\|_1 = \text{Tr}[\Delta_+] + \text{Tr}[\Delta_-] = 2\text{Tr}[\Delta_+].$$

So  $D(\rho, \sigma) = \text{Tr}[\Delta_+]$ .

Let  $P$  be the projector onto the support of  $\Delta_+$  with  $P\Delta_+ = \Delta_+$ . We evaluate the trace:

$$(3.110) \quad \text{Tr}[P\Delta] = \text{Tr}[P(\Delta_+ - \Delta_-)] = \text{Tr}[\Delta_+] = D(\rho, \sigma).$$

Therefore, the maximum is achieved.  $\square$

Proposition 3.38 implies that  $D(\rho, \sigma)$  is the maximum difference in the probability of obtaining a specific measurement outcome when measuring  $\rho$  versus  $\sigma$ .

A fundamental property of the trace distance is that it cannot increase under the action of a quantum channel. This reflects the physical intuition that noise or information loss (modeled by the channel) makes states harder to distinguish. This result parallels Proposition 3.28 for classical channels.

**THEOREM 3.39** (Quantum channels are contractive). *Let  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  be a quantum channel. For any  $\rho, \sigma \in \mathcal{D}(\mathbb{C}^N)$ ,*

$$(3.111) \quad D(\mathcal{Q}[\rho], \mathcal{Q}[\sigma]) \leq D(\rho, \sigma).$$

**PROOF.** Let  $\rho' = \mathcal{Q}[\rho]$  and  $\sigma' = \mathcal{Q}[\sigma]$ . By Proposition 3.38, there exists a projector  $P$  (specifically, onto the positive subspace of  $\rho' - \sigma'$ ) such that

$$(3.112) \quad D(\rho', \sigma') = \text{Tr}[P(\rho' - \sigma')] = \text{Tr}[P\mathcal{Q}[\rho - \sigma]].$$

Consider the decomposition  $\rho - \sigma = \Delta_+ - \Delta_-$ , where  $\Delta_+, \Delta_- \succeq 0$  are the positive and negative parts, respectively. As shown in the proof of Proposition 3.38,  $D(\rho, \sigma) = \text{Tr}[\Delta_+] = \text{Tr}[\Delta_-]$ .

Substituting the decomposition and using linearity:

$$(3.113) \quad D(\rho', \sigma') = \text{Tr}[P\mathcal{Q}[\Delta_+ - \Delta_-]] = \text{Tr}[P\mathcal{Q}[\Delta_+]] - \text{Tr}[P\mathcal{Q}[\Delta_-]].$$

We analyze the two terms. Since  $\mathcal{Q}$  is a positive map, and  $\Delta_- \succeq 0$ , the output  $\mathcal{Q}[\Delta_-]$  is positive semidefinite. Since  $P \succeq 0$ , the trace of the product of two positive operators is non-negative:  $\text{Tr}[P\mathcal{Q}[\Delta_-]] \geq 0$ . Therefore,

$$(3.114) \quad D(\rho', \sigma') \leq \text{Tr}[P\mathcal{Q}[\Delta_+]].$$

Next, since  $\mathcal{Q}[\Delta_+] \succeq 0$  and  $P \preceq I$ , we have  $I - P \succeq 0$ . Thus  $\text{Tr}[(I - P)\mathcal{Q}[\Delta_+]] \geq 0$ , which implies  $\text{Tr}[P\mathcal{Q}[\Delta_+]] \leq \text{Tr}[\mathcal{Q}[\Delta_+]]$ . Therefore,

$$(3.115) \quad D(\rho', \sigma') \leq \text{Tr}[\mathcal{Q}[\Delta_+]].$$

Finally, since  $\mathcal{Q}$  is trace-preserving,  $\text{Tr}[\mathcal{Q}[\Delta_+]] = \text{Tr}[\Delta_+]$ . Combining the inequalities, we obtain

$$(3.116) \quad D(\mathcal{Q}[\rho], \mathcal{Q}[\sigma]) \leq \text{Tr}[\Delta_+] = D(\rho, \sigma).$$

□

**3.5.3. Fidelity.** While the trace distance is an operationally useful metric for the distance between quantum states, another widely used measure is the fidelity. Fidelity quantifies the “overlap” between two quantum states, and generalizes the inner product between pure state vectors.

**Definition 3.40** (Fidelity). *The **fidelity** between two quantum states  $\rho, \sigma \in \mathcal{D}(\mathbb{C}^N)$  is defined as*

$$(3.117) \quad F(\rho, \sigma) := \text{Tr} \left[ \sqrt{\rho^{\frac{1}{2}} \sigma \rho^{\frac{1}{2}}} \right].$$

This definition can be rewritten using the trace norm. A more symmetric expression involves the operator  $A = \rho^{1/2} \sigma^{1/2}$ . Recall that the trace norm of  $A$  is  $\|A\|_1 = \text{Tr}[|A|] = \text{Tr}[\sqrt{A^\dagger A}]$ . Here  $A^\dagger A = \sigma^{1/2} \rho \sigma^{1/2}$ . The singular values of  $A$  are the square roots of the eigenvalues of  $A^\dagger A$  (and also  $AA^\dagger = \rho^{1/2} \sigma \rho^{1/2}$ ). Thus,

$$(3.118) \quad F(\rho, \sigma) = \left\| \rho^{1/2} \sigma^{1/2} \right\|_1.$$

This immediately establishes that fidelity is symmetric:  $F(\rho, \sigma) = F(\sigma, \rho)$ , since  $\|A\|_1 = \|A^\dagger\|_1$ .

**Remark 3.41.** Nomenclature can be confusing. Sometimes the quantity defined above is called the square root fidelity, and  $F(\rho, \sigma)^2$  is called the fidelity. The **infidelity** is then defined as  $1 - F(\rho, \sigma)^2$ . We will adhere to Definition 3.40. ◇

Fidelity satisfies  $0 \leq F(\rho, \sigma) \leq 1$ . The upper bound follows from Hölder's inequality (Proposition 3.31,  $p = q = 2$ ):

$$(3.119) \quad F(\rho, \sigma) = \left\| \rho^{1/2} \sigma^{1/2} \right\|_1 \leq \left\| \rho^{1/2} \right\|_2 \left\| \sigma^{1/2} \right\|_2.$$

Since  $\left\| \rho^{1/2} \right\|_2^2 = \text{Tr}[\rho^{1/2} \rho^{1/2}] = \text{Tr}[\rho] = 1$ , we have  $F(\rho, \sigma) \leq 1$ . Furthermore,  $F(\rho, \sigma) = 1$  if and only if  $\rho = \sigma$ .

Fidelity itself is not a distance metric (it does not satisfy the triangle inequality). However, it can be converted into a metric known as the angle or Bures angle.

**Definition 3.42** (Angle between quantum states). *The **angle** between two quantum states  $\rho, \sigma \in \mathcal{D}(\mathbb{C}^N)$  is*

$$(3.120) \quad \theta(\rho, \sigma) := \arccos(F(\rho, \sigma)) \in [0, \pi/2].$$

**Example 3.43** (Pure states). If  $\rho = |\psi\rangle\langle\psi|$  and  $\sigma = |\varphi\rangle\langle\varphi|$  are two pure states.

$$(3.121) \quad \rho^{1/2} \sigma \rho^{1/2} = |\psi\rangle\langle\psi| |\varphi\rangle\langle\varphi| |\psi\rangle\langle\psi| = |\langle\psi|\varphi\rangle|^2 |\psi\rangle\langle\psi|.$$

This is a rank-1 operator. Its only non-zero eigenvalue is  $|\langle\psi|\varphi\rangle|^2$ . The square root of this eigenvalue is  $|\langle\psi|\varphi\rangle|$ . Thus,

$$(3.122) \quad F(\rho, \sigma) = |\langle\psi|\varphi\rangle|.$$

The fidelity is the absolute value of the overlap between the state vectors.

More generally, if only one state is pure, say  $\rho = |\psi\rangle\langle\psi|$ , then

$$(3.123) \quad F(\rho, \sigma) = \sqrt{\langle\psi|\sigma|\psi\rangle}.$$

It is the square root of the overlap between the pure state  $|\psi\rangle$  and the mixed state  $\sigma$ .

Let us relate the trace distance and fidelity for pure states  $\rho, \sigma$ . Let the angle be  $\theta = \theta(\rho, \sigma)$ , so  $F(\rho, \sigma) = \cos \theta$ . We can choose a basis such that  $|\psi\rangle = |0\rangle$  and  $|\varphi\rangle = \cos \theta |0\rangle + \sin \theta |1\rangle$  (by adjusting global phase). In this 2D subspace, the difference  $\rho - \sigma$  is represented by the matrix:

$$(3.124) \quad \Delta = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} - \begin{pmatrix} \cos^2 \theta & \cos \theta \sin \theta \\ \cos \theta \sin \theta & \sin^2 \theta \end{pmatrix} = \begin{pmatrix} \sin^2 \theta & -\cos \theta \sin \theta \\ -\cos \theta \sin \theta & -\sin^2 \theta \end{pmatrix}.$$

The eigenvalues of  $\Delta$  are  $\pm \sin \theta$ . The trace norm is  $\|\Delta\|_1 = |\sin \theta| + |-\sin \theta| = 2 \sin \theta$  (since  $0 \leq \theta \leq \pi/2$ ).

$$(3.125) \quad D(\rho, \sigma) = \frac{1}{2} \|\rho - \sigma\|_1 = \sin \theta.$$

We can express this in terms of fidelity  $F = \cos \theta$ :

$$(3.126) \quad D(\rho, \sigma) = \sqrt{1 - F(\rho, \sigma)^2}.$$

◇

**Example 3.44** (Classical states). Let  $\rho, \sigma$  be classical states corresponding to probability distributions  $p, q$  (see Example 3.37). Since the operators are diagonal, the definition simplifies:

$$(3.127) \quad F(\rho, \sigma) = \sum_j \sqrt{p_j q_j}.$$



This is the classical Bhattacharyya coefficient. The relationship between trace distance and fidelity for classical states is characterized by the inequality:

$$\begin{aligned}
 D(\rho, \sigma) &= \frac{1}{2} \sum_j |p_j - q_j| \geq \frac{1}{2} \sum_j (\sqrt{p_j} - \sqrt{q_j})^2 \\
 &= \frac{1}{2} \sum_j (p_j + q_j - 2\sqrt{p_j q_j}) = 1 - \sum_j \sqrt{p_j q_j} = 1 - F(\rho, \sigma).
 \end{aligned}
 \tag{3.128}$$

The inequality step uses  $|a^2 - b^2| \geq (a - b)^2$  for  $a, b \geq 0$ .  $\diamond$

We have seen two extremes: for pure states  $D = \sqrt{1 - F^2}$ , while for classical states  $D \geq 1 - F$ . These relationships are generalized by the Fuchs–van de Graaf inequalities (see [NC00, Section 9.2]), which provide tight bounds relating the two measures for arbitrary quantum states.

**THEOREM 3.45** (Fuchs–van de Graaf inequalities). *For any  $\rho, \sigma \in \mathcal{D}(\mathbb{C}^N)$ ,*

$$1 - F(\rho, \sigma) \leq D(\rho, \sigma) \leq \sqrt{1 - F(\rho, \sigma)^2}.$$

We state a few important properties of fidelity without proof. Their proofs typically rely on a powerful result known as Uhlmann’s theorem, which relates the fidelity between two mixed states to the maximum overlap between their purifications (see [NC00, Chapter 9], [Wat18, Chapter 3]).

**Proposition 3.46** (Properties of Fidelity and Angle). *Let  $\rho, \sigma \in \mathcal{D}(\mathbb{C}^N)$ .*

- (1) *(Metric property) The angle  $\theta(\rho, \sigma)$  is a distance metric on  $\mathcal{D}(\mathbb{C}^N)$ .*
- (2) *(Contractivity) For any quantum channel  $\mathcal{Q}$ , the angle is contractive:*

$$\theta(\mathcal{Q}[\rho], \mathcal{Q}[\sigma]) \leq \theta(\rho, \sigma).$$

*Equivalently, fidelity increases (or stays the same) under quantum channels:*

$$F(\mathcal{Q}[\rho], \mathcal{Q}[\sigma]) \geq F(\rho, \sigma).$$

The Fuchs–van de Graaf inequalities (Theorem 3.45) can be rewritten in terms of the angle  $\theta = \theta(\rho, \sigma)$ :

$$2 \sin^2 \frac{\theta}{2} \leq D(\rho, \sigma) \leq \sin \theta.$$

When the states are close ( $\theta \ll 1$ ), we can use the approximations  $\sin \theta \approx \theta$  and  $2 \sin^2(\theta/2) \approx \theta^2/2$ . This gives

$$\frac{1}{2} \theta^2 \lesssim D(\rho, \sigma) \lesssim \theta.$$

This quadratic difference in scaling suggests that while the different distance metrics are related, they can behave very differently.

**Example 3.47.** Consider a target state  $\rho = |0\rangle\langle 0|$ . Let  $\theta \in [0, \pi/2]$  and define two pure states:

$$|\theta_+\rangle = \cos \theta |0\rangle + \sin \theta |1\rangle, \quad |\theta_-\rangle = \cos \theta |0\rangle - \sin \theta |1\rangle.$$

Let  $\sigma_+$  and  $\sigma_-$  be the corresponding density operators. We also consider the mixed state  $\sigma_M = \frac{1}{2}(\sigma_+ + \sigma_-)$ .

$$\sigma_M = \cos^2 \theta |0\rangle\langle 0| + \sin^2 \theta |1\rangle\langle 1|.$$

We compare the fidelities and trace distances to the target state  $\rho$ . The fidelities are identical:

$$F(\rho, \sigma_+) = F(\rho, \sigma_-) = F(\rho, \sigma_M) = \cos \theta.$$

However, the trace distances differ significantly. For the pure states (using Example 3.43):

$$(3.137) \quad D(\rho, \sigma_{\pm}) = \sin \theta.$$

For the mixed state  $\sigma_M$  (using Example 3.37):

$$(3.138) \quad D(\rho, \sigma_M) = \sin^2 \theta.$$

If  $\theta$  is small,  $D(\rho, \sigma_{\pm}) \approx \theta$  while  $D(\rho, \sigma_M) \approx \theta^2$ . The mixed state is quadratically closer to the target state in trace distance than its pure components, even though they all share the same fidelity. The coherent superpositions in  $\sigma_+$  and  $\sigma_-$  (the off-diagonal terms) cancel out in the incoherent mixture  $\sigma_M$ , leading to a state that is statistically closer to  $\rho$ .  $\diamond$

Which measure, fidelity or trace distance, is more physically relevant? The answer depends on the context. Fidelity can often be estimated experimentally (e.g., via the SWAP test), while estimating the trace distance generally requires full quantum state tomography.

On the other hand, the trace distance directly bounds the difference in measurement statistics. According to Proposition 3.38, the maximum difference in the probability of any measurement outcome  $M$  is bounded by the trace distance:

$$(3.139) \quad |\text{Tr}[M\rho] - \text{Tr}[M\sigma]| \leq D(\rho, \sigma).$$

If the trace distance is small, the states are statistically indistinguishable by any measurement.

### 3.6. Distance between quantum channels

Quantifying the distance between quantum channels is important for analyzing the precision of quantum gates, the robustness of quantum algorithms, and the distinguishability of physical processes. This section introduces the primary tools used for this purpose: the induced trace norm and the diamond norm.

**3.6.1. Induced trace norm.** We begin by considering norms induced on the space of linear maps (superoperators) by the Schatten norms on the input and output spaces.

**Definition 3.48.** For a linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$ , the **induced trace norm** (or the induced  $1 \rightarrow 1$  norm) is defined as

$$(3.140) \quad \|\mathcal{Q}\|_{1 \rightarrow 1} := \sup_{X \in L(\mathbb{C}^N), \|X\|_1 \leq 1} \|\mathcal{Q}[X]\|_1.$$

This norm quantifies the maximum amplification of the trace norm under the action of  $\mathcal{Q}$ .

Analogously, the **induced operator norm** (or the induced  $\infty \rightarrow \infty$  norm) is defined using the operator norm  $\|\cdot\|_{\infty}$ :

$$(3.141) \quad \|\mathcal{Q}\|_{\infty \rightarrow \infty} := \sup_{X \in L(\mathbb{C}^N), \|X\|_{\infty} \leq 1} \|\mathcal{Q}[X]\|_{\infty}.$$

Induced norms are inherently submultiplicative, a property useful when analyzing compositions of maps.

**Proposition 3.49** (Submultiplicativity). Let  $\mathcal{R} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^{N'})$  and  $\mathcal{Q} : L(\mathbb{C}^{N'}) \rightarrow L(\mathbb{C}^M)$  be linear maps. Then

$$(3.142) \quad \|\mathcal{Q} \circ \mathcal{R}\|_{1 \rightarrow 1} \leq \|\mathcal{Q}\|_{1 \rightarrow 1} \|\mathcal{R}\|_{1 \rightarrow 1}.$$

PROOF. For any  $X \in L(\mathbb{C}^N)$ , by the definition of the induced norm:

$$(3.143) \quad \|(\mathcal{Q} \circ \mathcal{R})(X)\|_1 = \|\mathcal{Q}[\mathcal{R}[X]]\|_1 \leq \|\mathcal{Q}\|_{1 \rightarrow 1} \|\mathcal{R}[X]\|_1 \leq \|\mathcal{Q}\|_{1 \rightarrow 1} \|\mathcal{R}\|_{1 \rightarrow 1} \|X\|_1.$$

Taking the supremum over  $X$  with  $\|X\|_1 \leq 1$  yields the result.  $\square$

To analyze these norms, we introduce the concept of the adjoint map. The space of linear operators  $L(\mathbb{C}^N)$  forms a Hilbert space under the Hilbert-Schmidt inner product  $\langle A, B \rangle = \text{Tr}(A^\dagger B)$ . The **adjoint map**  $\mathcal{Q}^\dagger : L(\mathbb{C}^M) \rightarrow L(\mathbb{C}^N)$  is uniquely defined by the relation

$$(3.144) \quad \langle Y, \mathcal{Q}(X) \rangle = \langle \mathcal{Q}^\dagger(Y), X \rangle,$$

for all  $X \in L(\mathbb{C}^N)$  and  $Y \in L(\mathbb{C}^M)$ .

The induced trace norm and the induced operator norm exhibit a duality relationship analogous to the duality between the trace norm and operator norm for matrices (Lemma 3.34).

**Proposition 3.50** (Duality of Induced Norms). *For any linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$ , the following duality relation holds:*

$$(3.145) \quad \|\mathcal{Q}\|_{1 \rightarrow 1} = \|\mathcal{Q}^\dagger\|_{\infty \rightarrow \infty}.$$

PROOF. We begin with the definition of the induced trace norm and apply the variational characterization of the trace norm (Lemma 3.34, Eq. (3.94)):

$$(3.146) \quad \begin{aligned} \|\mathcal{Q}\|_{1 \rightarrow 1} &= \sup_{\|X\|_1 \leq 1} \|\mathcal{Q}(X)\|_1 \\ &= \sup_{\|X\|_1 \leq 1} \left( \sup_{\|Y\|_\infty \leq 1} |\text{Tr}(Y^\dagger \mathcal{Q}[X])| \right). \end{aligned}$$

We exchange the order of the suprema and employ the definition of the adjoint map (Eq. (3.144)):

$$(3.147) \quad \|\mathcal{Q}\|_{1 \rightarrow 1} = \sup_{\|Y\|_\infty \leq 1} \left( \sup_{\|X\|_1 \leq 1} |\text{Tr}((\mathcal{Q}^\dagger(Y))^\dagger X)| \right).$$

The inner supremum is the characterization of the operator norm via duality (Lemma 3.34, Eq. (3.95)), applied to the operator  $W = \mathcal{Q}^\dagger(Y)$ . That is,  $\sup_{\|X\|_1 \leq 1} |\text{Tr}(W^\dagger X)| = \|W\|_\infty$ .

$$(3.148) \quad \|\mathcal{Q}\|_{1 \rightarrow 1} = \sup_{\|Y\|_\infty \leq 1} \|\mathcal{Q}^\dagger(Y)\|_\infty = \|\mathcal{Q}^\dagger\|_{\infty \rightarrow \infty}.$$

$\square$

To compute the induced trace norm, it is helpful to characterize the inputs that achieve the maximum. We first establish that for general linear maps, the maximum is attained on rank-1 operators.

**Lemma 3.51.** *For any linear map  $\mathcal{Q}$ , the induced  $1 \rightarrow 1$  norm is achieved by a rank-1 operator:*

$$(3.149) \quad \|\mathcal{Q}\|_{1 \rightarrow 1} = \sup\{\|\mathcal{Q}(|u\rangle\langle v|)\|_1 : \|u\|_2 = 1, \|v\|_2 = 1\}.$$

PROOF. Let  $C_1 = \{X : \|X\|_1 \leq 1\}$  be the unit ball in the trace norm. The function  $f(X) = \|\mathcal{Q}(X)\|_1$  is convex. Since  $C_1$  is a compact, convex set, the maximum of  $f(X)$  over  $C_1$  must be achieved at an extreme point of  $C_1$ . The extreme points of  $C_1$  are precisely the rank-1 operators of the form  $|u\rangle\langle v|$  with normalized vectors  $|u\rangle, |v\rangle$ .

Explicitly, let  $X$  maximize the norm, with  $\|X\|_1 = 1$ . Its SVD  $X = \sum_i s_i |u_i\rangle\langle v_i|$  is a convex combination (since  $\sum s_i = 1, s_i > 0$ ) of the rank-1 operators  $X_i = |u_i\rangle\langle v_i|$ . By the triangle inequality:

$$(3.150) \quad \|\mathcal{Q}(X)\|_1 = \left\| \sum_i s_i \mathcal{Q}(X_i) \right\|_1 \leq \sum_i s_i \|\mathcal{Q}(X_i)\|_1 \leq \max_i \|\mathcal{Q}(X_i)\|_1.$$

Thus, the maximum is achieved by one of the rank-1 operators  $X_i$ .  $\square$

We now investigate how these norms behave for positive maps. We first state the following result for positive maps without proof [Wat18, Eq. (3.329)].

**Lemma 3.52.** *Let  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  be a positive linear map. Then*

$$(3.151) \quad \|\mathcal{Q}\|_{1 \rightarrow 1} = \|\mathcal{Q}^\dagger(I_M)\|_\infty.$$

A celebrated result known as the Russo–Dye theorem [Wat18, Theorem 3.39] simplifies the calculation of the induced norm for such maps.

**THEOREM 3.53 (Russo–Dye).** *Let  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  be a positive linear map. Then*

$$(3.152) \quad \|\mathcal{Q}\|_{1 \rightarrow 1} = \max_{\|u\|_2=1} \text{Tr}(\mathcal{Q}(|u\rangle\langle u|)).$$

**PROOF.** Since  $\mathcal{Q}^\dagger(I_M)$  is Hermitian (in fact positive semidefinite), its operator norm is the largest eigenvalue:

$$(3.153) \quad \|\mathcal{Q}^\dagger(I_M)\|_\infty = \sup_{\|u\|_2=1} \langle u | \mathcal{Q}^\dagger(I_M) | u \rangle = \sup_{\|u\|_2=1} \text{Tr}(\mathcal{Q}(|u\rangle\langle u|)).$$

The result follows from Lemma 3.52.  $\square$

As an immediate consequence, if  $\mathcal{Q}$  is a quantum channel, it is positive and trace-preserving. Thus,

$$(3.154) \quad \|\mathcal{Q}\|_{1 \rightarrow 1} = \max_u \text{Tr}(\mathcal{Q}(|u\rangle\langle u|)) = \max_u \text{Tr}(|u\rangle\langle u|) = 1.$$

The fact that quantum channels have an induced trace norm of 1 leads to an important stability property for compositions of channels.

**Proposition 3.54.** *Let  $\mathcal{Q}_1, \dots, \mathcal{Q}_K$  and  $\tilde{\mathcal{Q}}_1, \dots, \tilde{\mathcal{Q}}_K$  be sequences of quantum channels. Then*

$$(3.155) \quad \left\| \mathcal{Q}_K \circ \dots \circ \mathcal{Q}_1 - \tilde{\mathcal{Q}}_K \circ \dots \circ \tilde{\mathcal{Q}}_1 \right\|_{1 \rightarrow 1} \leq \sum_{i=1}^K \left\| \mathcal{Q}_i - \tilde{\mathcal{Q}}_i \right\|_{1 \rightarrow 1}.$$

**PROOF.** We use a telescoping sum argument. For  $K = 2$ :

$$(3.156) \quad \mathcal{Q}_2 \circ \mathcal{Q}_1 - \tilde{\mathcal{Q}}_2 \circ \tilde{\mathcal{Q}}_1 = (\mathcal{Q}_2 - \tilde{\mathcal{Q}}_2) \circ \mathcal{Q}_1 + \tilde{\mathcal{Q}}_2 \circ (\mathcal{Q}_1 - \tilde{\mathcal{Q}}_1).$$

By the triangle inequality and submultiplicativity (Proposition 3.49):

$$(3.157) \quad \begin{aligned} \left\| \mathcal{Q}_2 \circ \mathcal{Q}_1 - \tilde{\mathcal{Q}}_2 \circ \tilde{\mathcal{Q}}_1 \right\|_{1 \rightarrow 1} &\leq \left\| \mathcal{Q}_2 - \tilde{\mathcal{Q}}_2 \right\|_{1 \rightarrow 1} \|\mathcal{Q}_1\|_{1 \rightarrow 1} \\ &\quad + \left\| \tilde{\mathcal{Q}}_2 \right\|_{1 \rightarrow 1} \left\| \mathcal{Q}_1 - \tilde{\mathcal{Q}}_1 \right\|_{1 \rightarrow 1}. \end{aligned}$$

Since  $\mathcal{Q}_1$  and  $\tilde{\mathcal{Q}}_2$  are quantum channels, their induced trace norms are 1.

$$(3.158) \quad \left\| \mathcal{Q}_2 \circ \mathcal{Q}_1 - \tilde{\mathcal{Q}}_2 \circ \tilde{\mathcal{Q}}_1 \right\|_{1 \rightarrow 1} \leq \left\| \mathcal{Q}_2 - \tilde{\mathcal{Q}}_2 \right\|_{1 \rightarrow 1} + \left\| \mathcal{Q}_1 - \tilde{\mathcal{Q}}_1 \right\|_{1 \rightarrow 1}.$$

The general case follows by induction.  $\square$

**3.6.2. The diamond norm.** The induced trace norm quantifies how much a map  $\mathcal{Q}$  acting on a system  $S$  changes the state of  $S$ . However, this is insufficient in quantum mechanics due to entanglement. If  $S$  is entangled with an auxiliary system  $A$ , the action of  $\mathcal{Q}$  on  $S$  (described by  $\mathcal{Q} \otimes \mathcal{I}_A$ ) might alter the joint state of  $SA$  significantly more than predicted by  $\|\mathcal{Q}\|_{1 \rightarrow 1}$ . To capture the true behavior of the map in the presence of arbitrary entanglement, we must consider its stabilized action. This leads to the **diamond norm**, also known as the **completely bounded trace norm**.

**Definition 3.55** (Diamond Norm). *Let  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$  be a linear map. The diamond norm of  $\mathcal{Q}$  is defined as*

$$(3.159) \quad \|\mathcal{Q}\|_{\diamond} := \sup_{k \geq 1} \|\mathcal{Q} \otimes \mathcal{I}_k\|_{1 \rightarrow 1} = \sup_{k \geq 1} \|\mathcal{I}_k \otimes \mathcal{Q}\|_{1 \rightarrow 1},$$

where  $\mathcal{I}_k$  denotes the identity map on  $L(\mathbb{C}^k)$ .

If  $\mathcal{Q}$  is a quantum channel, then for every  $k$  the map  $\mathcal{Q} \otimes \mathcal{I}_k$  is also a quantum channel, and hence has induced trace norm 1. Therefore,

$$(3.160) \quad \|\mathcal{Q}\|_{\diamond} = \sup_{k \geq 1} \|\mathcal{Q} \otimes \mathcal{I}_k\|_{1 \rightarrow 1} = 1.$$

While the definition involves a supremum over all dimensions  $k$ , a remarkable result shows that the supremum is achieved when the auxiliary dimension matches the input dimension of the map.

**Proposition 3.56** (Stabilization of the Diamond Norm). *For any linear map  $\mathcal{Q} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^M)$ , the supremum in Eq. (3.159) is achieved for  $k = N$ . That is,*

$$(3.161) \quad \|\mathcal{Q}\|_{\diamond} = \|\mathcal{Q} \otimes \mathcal{I}_N\|_{1 \rightarrow 1}.$$

PROOF. We aim to show that for any  $k \geq 1$ ,  $\|\mathcal{Q} \otimes \mathcal{I}_k\|_{1 \rightarrow 1} \leq \|\mathcal{Q} \otimes \mathcal{I}_N\|_{1 \rightarrow 1}$ .

Let  $k \geq 1$ . By Lemma 3.51, the induced norm is achieved by a rank-1 input. There exist normalized vectors  $|\alpha\rangle, |\beta\rangle \in \mathbb{C}^N \otimes \mathbb{C}^k$  such that

$$(3.162) \quad \|\mathcal{Q} \otimes \mathcal{I}_k\|_{1 \rightarrow 1} = \|(\mathcal{Q} \otimes \mathcal{I}_k)(|\alpha\rangle\langle\beta|)\|_1.$$

Consider the Schmidt decompositions of  $|\alpha\rangle$  and  $|\beta\rangle$ . The Schmidt ranks  $r, s$  are at most  $N$ .

$$(3.163) \quad |\alpha\rangle = \sum_{i=1}^r \sqrt{p_i} |a_i\rangle \otimes |x_i\rangle, \quad |\beta\rangle = \sum_{j=1}^s \sqrt{q_j} |b_j\rangle \otimes |y_j\rangle.$$

Here,  $\{|a_i\rangle\}, \{|b_j\rangle\} \subset \mathbb{C}^N$  and  $\{|x_i\rangle\}, \{|y_j\rangle\} \subset \mathbb{C}^k$  are orthonormal sets. Let  $Y = (\mathcal{Q} \otimes \mathcal{I}_k)(|\alpha\rangle\langle\beta|)$ .

$$(3.164) \quad Y = \sum_{i,j} \sqrt{p_i q_j} \mathcal{Q}(|a_i\rangle\langle b_j|) \otimes |x_i\rangle\langle y_j|.$$

We construct corresponding vectors in  $\mathbb{C}^N \otimes \mathbb{C}^N$ . Let  $\{|e_i\rangle\}_{i=1}^N$  be a basis for  $\mathbb{C}^N$ . Define normalized vectors  $|\alpha'\rangle, |\beta'\rangle \in \mathbb{C}^N \otimes \mathbb{C}^N$  by replacing  $|x_i\rangle$  with  $|e_i\rangle$  and  $|y_j\rangle$  with  $|e_j\rangle$ . Let  $Y' = (\mathcal{Q} \otimes \mathcal{I}_N)(|\alpha'\rangle\langle\beta'|)$ .

$$(3.165) \quad Y' = \sum_{i,j} \sqrt{p_i q_j} \mathcal{Q}(|a_i\rangle\langle b_j|) \otimes |e_i\rangle\langle e_j|.$$

We show that  $\|Y\|_1 = \|Y'\|_1$ . Define partial isometries  $V, W : \mathbb{C}^N \rightarrow \mathbb{C}^k$ . Let  $V$  map  $\text{span}\{|e_i\rangle\}_{i=1}^r$  isometrically onto  $\text{span}\{|x_i\rangle\}_{i=1}^r$ , and similarly for  $W$  and  $\{|y_j\rangle\}$ . We can relate  $Y$  and  $Y'$ :

$$(3.166) \quad Y = (I_M \otimes V)Y'(I_M \otimes W^\dagger).$$

Extend  $V$  and  $W$  to unitaries  $\tilde{V}, \tilde{W}$  on  $\mathbb{C}^k$  (by choosing orthonormal complements). Since  $Y'$  only has support on the subspaces where  $V$  and  $W$  act isometrically, we have

$$(3.167) \quad Y = (I_M \otimes \tilde{V})Y'(I_M \otimes \tilde{W}^\dagger).$$

By unitary invariance of the trace norm,  $\|Y\|_1 = \|Y'\|_1$ .

We have established  $\|\mathcal{Q} \otimes \mathcal{I}_k\|_{1 \rightarrow 1} = \|Y'\|_1$ . Since  $\|\alpha'\langle\beta'\|_1 = 1$ , we have  $\|Y'\|_1 \leq \|\mathcal{Q} \otimes \mathcal{I}_N\|_{1 \rightarrow 1}$ . This completes the proof.  $\square$

The diamond norm inherits the submultiplicativity property from the induced trace norm.

**Proposition 3.57** (Submultiplicativity of the Diamond Norm). *Let  $\mathcal{R} : L(\mathbb{C}^N) \rightarrow L(\mathbb{C}^{N'})$  and  $\mathcal{Q} : L(\mathbb{C}^{N'}) \rightarrow L(\mathbb{C}^M)$  be linear maps. Then*

$$(3.168) \quad \|\mathcal{Q} \circ \mathcal{R}\|_\diamond \leq \|\mathcal{Q}\|_\diamond \|\mathcal{R}\|_\diamond.$$

PROOF. We use the definition of the diamond norm and the property that  $(\mathcal{Q} \circ \mathcal{R}) \otimes \mathcal{I}_k = (\mathcal{Q} \otimes \mathcal{I}_k) \circ (\mathcal{R} \otimes \mathcal{I}_k)$ .

$$(3.169) \quad \|\mathcal{Q} \circ \mathcal{R}\|_\diamond = \sup_k \|(\mathcal{Q} \otimes \mathcal{I}_k) \circ (\mathcal{R} \otimes \mathcal{I}_k)\|_{1 \rightarrow 1}.$$

By the submultiplicativity of the induced trace norm (Proposition 3.49):

$$(3.170) \quad \begin{aligned} \|\mathcal{Q} \circ \mathcal{R}\|_\diamond &\leq \sup_k (\|\mathcal{Q} \otimes \mathcal{I}_k\|_{1 \rightarrow 1} \|\mathcal{R} \otimes \mathcal{I}_k\|_{1 \rightarrow 1}) \\ &\leq \left( \sup_k \|\mathcal{Q} \otimes \mathcal{I}_k\|_{1 \rightarrow 1} \right) \left( \sup_k \|\mathcal{R} \otimes \mathcal{I}_k\|_{1 \rightarrow 1} \right) \\ &= \|\mathcal{Q}\|_\diamond \|\mathcal{R}\|_\diamond. \end{aligned}$$

$\square$

We can derive useful bounds on the diamond norm for specific types of maps. We start with maps defined by a single Kraus operator.

**Lemma 3.58.** *Let  $\mathcal{Q}_A(X) = AXA^\dagger$  and  $\mathcal{Q}_B(X) = BXB^\dagger$ . Then the diamond norm of their difference is bounded by*

$$(3.171) \quad \|\mathcal{Q}_A - \mathcal{Q}_B\|_\diamond \leq (\|A\|_\infty + \|B\|_\infty) \|A - B\|_\infty.$$

PROOF. Let  $\Phi = \mathcal{Q}_A - \mathcal{Q}_B$ . By stabilization (Proposition 3.56), we evaluate  $\|\Phi \otimes \mathcal{I}_N\|_{1 \rightarrow 1}$ . Let  $X_{SR}$  be an input operator with  $\|X_{SR}\|_1 = 1$ .

$$(3.172) \quad (\Phi \otimes \mathcal{I}_N)(X_{SR}) = (A \otimes I)X_{SR}(A^\dagger \otimes I) - (B \otimes I)X_{SR}(B^\dagger \otimes I).$$

We use the identity  $AA^\dagger - BB^\dagger = A(A^\dagger - B^\dagger) + (A - B)B^\dagger$ .

$$(3.173) \quad \begin{aligned} (\Phi \otimes \mathcal{I}_N)(X_{SR}) &= (A \otimes I)X_{SR}((A^\dagger - B^\dagger) \otimes I) \\ &\quad + ((A - B) \otimes I)X_{SR}(B^\dagger \otimes I). \end{aligned}$$

We bound the trace norm using the triangle inequality and Hölder's inequality ( $\|Y_1XY_2\|_1 \leq \|Y_1\|_\infty \|X\|_1 \|Y_2\|_\infty$ ). Since  $\|X_{SR}\|_1 = 1$  and  $\|Y \otimes I\|_\infty = \|Y\|_\infty$ :

$$(3.174) \quad \|(\Phi \otimes \mathcal{I}_N)(X_{SR})\|_1 \leq \|A\|_\infty \|A^\dagger - B^\dagger\|_\infty + \|A - B\|_\infty \|B^\dagger\|_\infty.$$

Using  $\|A^\dagger - B^\dagger\|_\infty = \|A - B\|_\infty$  and  $\|B^\dagger\|_\infty = \|B\|_\infty$ , we obtain the bound.  $\square$

**Example 3.59** (Distance between unitary channels). Consider unitary channels  $\mathcal{U}(X) = UXU^\dagger$  and  $\mathcal{V}(X) = VXV^\dagger$ . Since  $\|U\|_\infty = \|V\|_\infty = 1$ , Lemma 3.58 yields the bound:

$$(3.175) \quad \|\mathcal{U} - \mathcal{V}\|_\diamond \leq 2\|U - V\|_\infty.$$

$\diamond$

While the bound in Eq. (3.175) is widely used, it is not always tight. Furthermore, one might expect that stabilization is necessary for unitary channels. However, the difference between unitary channels exhibits a special structure that renders stabilization unnecessary.

**Proposition 3.60.** *Let  $\mathcal{U}, \mathcal{V}$  be two unitary channels defined by unitaries  $U$  and  $V$ . Then the diamond norm of their difference is equal to the induced trace norm:*

$$(3.176) \quad \|\mathcal{U} - \mathcal{V}\|_\diamond = \|\mathcal{U} - \mathcal{V}\|_{1 \rightarrow 1}.$$

*This norm can be computed explicitly using the numerical range of  $W = U^\dagger V$ :*

$$(3.177) \quad \|\mathcal{U} - \mathcal{V}\|_\diamond = 2\sqrt{1 - d_{\min}^2},$$

where  $d_{\min} = \inf\{|z| : z \in \mathcal{W}(W)\}$  is the minimum distance from the origin to the numerical range  $\mathcal{W}(W) = \{\langle x|W|x\rangle : \|x\|_2 = 1\}$ .

**PROOF.** Let  $\Phi = \mathcal{U} - \mathcal{V}$ . We first establish a lower bound for the induced trace norm  $\|\Phi\|_{1 \rightarrow 1}$ . According to Lemma 3.51, the induced trace norm is defined by the supremum over rank-1 inputs. Restricting the optimization to pure states  $\rho = |x\rangle\langle x|$  yields a lower bound:

$$(3.178) \quad \|\Phi\|_{1 \rightarrow 1} \geq \sup_{|x\rangle} \|\Phi(|x\rangle\langle x|)\|_1.$$

The output is

$$(3.179) \quad \Phi(|x\rangle\langle x|) = U|x\rangle\langle x|U^\dagger - V|x\rangle\langle x|V^\dagger = |\psi_U\rangle\langle\psi_U| - |\psi_V\rangle\langle\psi_V|,$$

where  $|\psi_U\rangle = U|x\rangle$  and  $|\psi_V\rangle = V|x\rangle$ . The trace norm of the difference between two pure states is determined by their overlap (see Example 3.43):

$$(3.180) \quad \| |\psi_U\rangle\langle\psi_U| - |\psi_V\rangle\langle\psi_V| \|_1 = 2\sqrt{1 - |\langle\psi_U|\psi_V\rangle|^2}.$$

The overlap is  $\langle\psi_U|\psi_V\rangle = \langle x|U^\dagger V|x\rangle = \langle x|W|x\rangle$ . To maximize the norm, we must minimize the magnitude of the overlap. The set of values  $\{\langle x|W|x\rangle : \|x\|_2 = 1\}$  is the numerical range  $\mathcal{W}(W)$ . Thus, the supremum over pure states is

$$(3.181) \quad 2\sqrt{1 - \inf_{z \in \mathcal{W}(W)} |z|^2} = 2\sqrt{1 - d_{\min}^2}.$$

Next, we consider the diamond norm  $\|\Phi\|_\diamond$ . By the stabilization property (Proposition 3.56),  $\|\Phi\|_\diamond = \|\Phi \otimes \mathcal{I}_N\|_{1 \rightarrow 1}$ . Unlike the induced trace norm, the diamond norm is achieved on pure states (see [Wat18, Theorem 3.51]). Let  $|\Psi\rangle \in \mathbb{C}^N \otimes \mathbb{C}^N$  be a normalized pure state. The action of the map

on  $\rho = |\Psi\rangle\langle\Psi|$  yields the difference of two pure states  $|\Psi_U\rangle = (U \otimes I)|\Psi\rangle$  and  $|\Psi_V\rangle = (V \otimes I)|\Psi\rangle$ . The norm is again given by  $2\sqrt{1 - |\langle\Psi_U|\Psi_V\rangle|^2}$ . The overlap is

$$(3.182) \quad \langle\Psi_U|\Psi_V\rangle = \langle\Psi|(U^\dagger \otimes I)(V \otimes I)|\Psi\rangle = \langle\Psi|(W \otimes I)|\Psi\rangle.$$

We express this overlap in terms of the reduced density operator  $\rho_A = \text{Tr}_B[|\Psi\rangle\langle\Psi|]$ :

$$(3.183) \quad \langle\Psi|(W \otimes I)|\Psi\rangle = \text{Tr}[(W \otimes I)|\Psi\rangle\langle\Psi|] = \text{Tr}[W\rho_A].$$

As  $|\Psi\rangle$  varies over all pure states in the joint space,  $\rho_A$  varies over all density operators in  $\mathcal{D}(\mathbb{C}^N)$ . The set of achievable overlaps is therefore the set of expectation values  $\{\text{Tr}[W\rho] : \rho \in \mathcal{D}(\mathbb{C}^N)\}$ . This set is the convex hull of the numerical range  $\mathcal{W}(W)$ . By the Toeplitz–Hausdorff theorem (see [Bha97, Chapter 1]), the numerical range  $\mathcal{W}(W)$  is a convex set. Therefore, the convex hull of  $\mathcal{W}(W)$  is  $\mathcal{W}(W)$  itself. This implies that allowing entanglement does not extend the range of possible overlaps:

$$(3.184) \quad \inf_{\|\Psi\|_2=1} |\langle\Psi|(W \otimes I)|\Psi\rangle| = \inf_{z \in \mathcal{W}(W)} |z| = d_{\min}.$$

Consequently,

$$(3.185) \quad \|\Phi\|_\diamond = 2\sqrt{1 - d_{\min}^2}.$$

Combining this with Eq. (3.181) and the inequality  $\|\Phi\|_{1 \rightarrow 1} \leq \|\Phi\|_\diamond$ , we conclude  $\|\Phi\|_\diamond = \|\Phi\|_{1 \rightarrow 1}$ .  $\square$

**Example 3.61.** Consider the  $2 \times 2$  unitaries  $U = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$  and  $V = I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ . We calculate the operator norm of their difference:

$$(3.186) \quad U - V = \begin{pmatrix} -1 & 1 \\ -1 & -1 \end{pmatrix}.$$

The singular values are the square roots of the eigenvalues of  $(U - V)^\dagger(U - V) = \text{diag}(2, 2)$ . Thus,  $\|U - V\|_\infty = \sqrt{2}$ . The general bound in Eq. (3.175) gives  $\|\mathcal{U} - \mathcal{V}\|_\diamond \leq 2\sqrt{2} \approx 2.828$ .

However, as  $W = U^\dagger = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$ , the eigenvalues of  $W$  are  $i$  and  $-i$ . Since  $W$  is normal,  $\mathcal{W}(W)$  is the convex hull of the eigenvalues, i.e., the segment  $[-i, i]$  on the imaginary axis. The minimum distance to the origin is  $d_{\min} = 0$ . Thus, the exact diamond norm is  $2\sqrt{1 - 0^2} = 2$ .  $\diamond$

**Example 3.62** (Qubit Phase Shift Channel). We illustrate the computation using a single-qubit example. Consider the identity channel  $\mathcal{I}$  ( $U = I$ ) and the phase shift channel  $\mathcal{P}_\theta$ , defined by the unitary  $V = P_\theta = \text{diag}(1, e^{i\theta})$ . We wish to compute  $\|\mathcal{I} - \mathcal{P}_\theta\|_\diamond$ .

We apply Proposition 3.60. We compute  $W = U^\dagger V = P_\theta$ . We need to determine the numerical range  $\mathcal{W}(P_\theta)$ . Since  $P_\theta$  is a normal operator, its numerical range is the convex hull of its eigenvalues,  $\{1, e^{i\theta}\}$ . This is the line segment (chord) connecting 1 and  $e^{i\theta}$  in the complex plane.

We seek the minimum distance  $d_{\min}$  from the origin to this segment. Geometrically, this distance is the altitude of the isosceles triangle formed by the origin and the two eigenvalues.

The length of the base of the triangle (the chord) is  $|1 - e^{i\theta}| = \sqrt{(1 - \cos\theta)^2 + \sin^2\theta} = \sqrt{2 - 2\cos\theta} = 2|\sin(\theta/2)|$ . The area of the triangle is  $\frac{1}{2}|\sin\theta|$ . Let  $h$  be the altitude, which corresponds to  $d_{\min}$ . The area is also  $\frac{1}{2} \cdot \text{base} \cdot h$ .

$$(3.187) \quad d_{\min} = h = \frac{|\sin\theta|}{2|\sin(\theta/2)|} = \frac{2|\sin(\theta/2)\cos(\theta/2)|}{2|\sin(\theta/2)|} = |\cos(\theta/2)|.$$



Substituting this minimum value into Eq. (3.177):

$$(3.188) \quad \|\mathcal{I} - \mathcal{P}_\theta\|_\diamond = 2\sqrt{1 - \cos^2(\theta/2)} = 2\sqrt{\sin^2(\theta/2)} = 2|\sin(\theta/2)|.$$

By Proposition 3.60, the induced trace norm is identical:  $\|\mathcal{I} - \mathcal{P}_\theta\|_{1 \rightarrow 1} = 2|\sin(\theta/2)|$ .

For instance, if  $\theta = \pi$ , the channel is the Pauli-Z channel  $\mathcal{Z}$ . The diamond norm is  $2|\sin(\pi/2)| = 2$ . The minimum overlap is  $d_{\min} = 0$ . This is achieved by the input state  $|+\rangle = \frac{1}{\sqrt{2}}(|0\rangle + |1\rangle)$ , since  $\langle + | \mathcal{Z} | + \rangle = 0$ .  $\diamond$

The following example illustrates that the standard induced trace norm can drastically underestimate the “size” of a map that is not completely positive.

**Example 3.63** (Transpose Map). Let  $\mathcal{T} : \mathbb{C}^{N \times N} \rightarrow \mathbb{C}^{N \times N}$  be the transpose map,  $\mathcal{T}(X) = X^\top$ , defined in a fixed basis. Since the transpose preserves the eigenvalues of Hermitian matrices and maps density matrices to density matrices, it preserves the 1-norm for positive inputs. It can be shown that  $\|\mathcal{T}\|_{1 \rightarrow 1} = 1$ .

However, consider the action of  $\mathcal{T} \otimes \mathcal{I}_N$  on the unnormalized maximally entangled state  $|\Omega\rangle = \sum_{i=1}^N |i\rangle \otimes |i\rangle$ . The corresponding density matrix is  $\omega = \sum_{i,j} |i\rangle\langle j| \otimes |i\rangle\langle j|$ . Applying the partial transpose yields

$$(3.189) \quad (\mathcal{T} \otimes \mathcal{I}_N)(\omega) = \sum_{i,j} |j\rangle\langle i| \otimes |i\rangle\langle j|,$$

which is the SWAP operator. The eigenvalues of the SWAP operator on  $\mathbb{C}^N \otimes \mathbb{C}^N$  are  $+1$  (on the symmetric subspace of dimension  $N(N+1)/2$ ) and  $-1$  (on the antisymmetric subspace of dimension  $N(N-1)/2$ ). The trace norm is the sum of singular values (absolute values of eigenvalues):

$$(3.190) \quad \|(\mathcal{T} \otimes \mathcal{I}_N)(\omega)\|_1 = \frac{N(N+1)}{2} + \frac{N(N-1)}{2} = N^2.$$

Since  $\|\omega\|_1 = \|\Omega\|^2 = N$ , we find that for this specific state, the ratio of output norm to input norm is  $N$ . Thus  $\|\mathcal{T}\|_\diamond \geq N$ .  $\diamond$

**3.6.3. Induced trace distance and diamond distance.** The induced trace distance between two linear maps  $\mathcal{Q}, \mathcal{R}$  is

$$(3.191) \quad D(\mathcal{Q}, \mathcal{R}) := \frac{1}{2} \|\mathcal{Q} - \mathcal{R}\|_{1 \rightarrow 1}.$$

**Example 3.64** (Trace distance for classical channel). Given two transition matrices  $Q, R \in \mathbb{R}^{N \times N}$ , the corresponding classical channels are

$$(3.192) \quad \mathcal{Q}[\rho] = \sum_{i,j \in [N]} Q_{ij} |i\rangle\langle j| \rho |j\rangle\langle i|, \quad \mathcal{R}[\rho] = \sum_{i,j \in [N]} R_{ij} |i\rangle\langle j| \rho |j\rangle\langle i|.$$

Then

$$\begin{aligned}
D(\mathcal{Q}, \mathcal{R}) &= \frac{1}{2} \sup_{\|\rho\|_1=1} \|\mathcal{Q}[\rho] - \mathcal{R}[\rho]\|_1 \\
&= \frac{1}{2} \sup_{\|\rho\|_1=1} \sum_i \left| \sum_j (Q_{ij} - R_{ij}) \rho_{jj} \right| \\
(3.193) \quad &\leq \frac{1}{2} \sup_{\|\rho\|_1=1} \left( \max_j \sum_i |Q_{ij} - R_{ij}| \right) \text{Tr} |\rho| \\
&\leq \frac{1}{2} \sup_{\|\rho\|_1=1} \left( \max_j \sum_i |Q_{ij} - R_{ij}| \right) \|\rho\|_1 \\
&= D(Q, R),
\end{aligned}$$

which is the induced total variation distance between the transition matrices  $Q, R$ . Here we have used Proposition 3.30 in the last inequality. On the other hand, choosing  $\rho = |j'\rangle\langle j'|$  with  $j' = \arg \max_j \sum_i |Q_{ij} - R_{ij}|$ , we have  $D(\mathcal{Q}, \mathcal{R}) \geq D(Q, R)$ . This proves that the induced trace distance is consistent with the induced total variation distance on classical channels:

$$(3.194) \quad D(\mathcal{Q}, \mathcal{R}) = D(Q, R).$$

◇

The metric induced by the diamond norm is known as the diamond distance. The factor of  $1/2$  normalizes the metric such that perfectly distinguishable channels have a distance of 1, analogous to the trace distance for quantum states.

**Definition 3.65** (Diamond Distance). *Let  $\mathcal{Q}, \mathcal{R} : \mathbb{C}^{N \times N} \rightarrow \mathbb{C}^{M \times M}$  be two linear maps. The **diamond distance** between them is defined as*

$$(3.195) \quad D_\diamond(\mathcal{Q}, \mathcal{R}) := \frac{1}{2} \|\mathcal{Q} - \mathcal{R}\|_\diamond.$$

Quantum channels satisfy the linear error growth property with respect to the diamond distance. The proof is also very similar to Proposition 3.54.

**Proposition 3.66.** *Let  $\{\mathcal{U}_i\}_{i=1}^K$  and  $\{\tilde{\mathcal{U}}_i\}_{i=1}^K$  be sequences of unitary channels generated by the unitary operators  $\{U_i\}_{i=1}^K$  and  $\{\tilde{U}_i\}_{i=1}^K$ , respectively. The diamond distance between the composite channels is bounded by*

$$(3.196) \quad D_\diamond(\mathcal{U}_K \cdots \mathcal{U}_1, \tilde{\mathcal{U}}_K \cdots \tilde{\mathcal{U}}_1) \leq \sum_{i=1}^K \left\| U_i - \tilde{U}_i \right\|_\infty.$$

PROOF. First, we observe that quantum channels satisfy a linear error growth property with respect to the diamond distance. The proof of this property relies on a telescoping sum argument, which is strictly analogous to the proof of Proposition 3.54 and is therefore omitted. This yields the bound

$$(3.197) \quad D_\diamond(\mathcal{U}_K \cdots \mathcal{U}_1, \tilde{\mathcal{U}}_K \cdots \tilde{\mathcal{U}}_1) \leq \frac{1}{2} \sum_{i=1}^K \left\| \tilde{\mathcal{U}}_i - \mathcal{U}_i \right\|_\diamond.$$

It suffices to bound the diamond norm difference for a single step. Recalling Eq. (3.175), we have the general bound  $\|\tilde{\mathcal{U}}_i - \mathcal{U}_i\|_{\diamond} \leq 2 \|U_i - \tilde{U}_i\|$ . Substituting this estimate into the linear error growth inequality completes the proof.  $\square$

### Notes and further reading

The formalism of quantum channels rests on foundational results in operator theory. The operator-sum representation (Theorem 3.18) is due to Kraus [KBDW83], while the dilation theorem (Theorem 3.20) was established by Stinespring [Sti55]. The isomorphism characterizing completely positive maps via their action on entangled states is attributed to Choi [Cho75] and Jamiołkowski [Jam72].

The induced trace distance provides a useful way to compare two channels via their action on input states. It is worth noting that the contractivity properties of the trace distance used in this context rely on positivity and trace preservation, and do not require complete positivity. By contrast, complete positivity is required to ensure that a channel remains positive when extended by an identity map on an arbitrary ancillary register. This distinction becomes operationally visible in the channel discrimination task: for some pairs of channels, optimal discrimination is only possible when the input is entangled with an ancillary register. This motivates the use of stabilized distances such as the diamond norm (the completely bounded trace norm), which explicitly accounts for ancillary extensions. For distance measures, Helstrom [Hel69] provided the operational interpretation of the trace distance in terms of state discrimination. Fidelity was studied by Uhlmann [Uhl76] as transition probability. The tight relationship between these two measures (Theorem 3.45) was established by Fuchs and van de Graaf [FVDG02]. The diamond norm was introduced to quantum computing by Kitaev [Kit97] to quantify the accuracy of quantum gates in a manner robust to entanglement, and is closely related to the completely bounded norm in operator algebra. We refer readers to [Wat18, Chapter 3.3] for further discussion.

Most of the discussions in this book will be restricted to unitary channels, and these unitary channels are often applied to pure states. Nevertheless, the concept of a quantum channel is helpful for understanding the probabilistic nature of quantum algorithms. For a systematic treatment of density operators and quantum channels, we refer readers to [Wat18, Chapter 2] and [NC00, Section 2.4, 8.2]. We refer readers to [Wat18, Chapter 3] for properties of the norms and distances introduced here, and their applications in discrimination-type problems. For matrix analysis tools, such as Schatten norms, we refer to [Bha97].



## CHAPTER 4

# Universality of quantum circuits



## Quantum processing of classical information

Quantum algorithms often require classical data to be loaded, processed, and manipulated within a quantum circuit. This chapter explores how classical information can be encoded and operated on in a quantum computing framework. We begin with the reversible simulation of classical logic gates, a prerequisite for embedding classical computation into quantum circuits. We then discuss uncomputation, which is very useful for cleaning up intermediate states without disturbing the computation's outcome. The chapter proceeds to cover fixed-point number representation and quantum random access memory (QRAM). Finally, we present methods for implementing certain classical arithmetic operations within quantum circuits.

### 5.1. Reversible simulation of classical gates

How can we compare the computational power of quantum computers to that of classical computers? While it remains extremely difficult to prove that quantum computers are fundamentally more powerful than classical ones, it is well established that quantum computers are at least as powerful. More precisely, any classical circuit can be simulated asymptotically efficiently by a quantum circuit.

The key idea behind this equivalence lies in the reversible simulation of classical gates. Some classical logic gates, such as the NOT gate, are already reversible and can be directly implemented by the Pauli X gate. However, many commonly used gates, including AND, OR, and NAND, are not reversible and cannot be directly translated into unitary transformations. This leads to a foundational step: expressing classical computation in terms of reversible logic gates.

**Reversible computation**, which predates quantum computing, was originally studied in the context of thermodynamics and the fundamental limits of energy dissipation. It refers to models of computation in which each operation can be uniquely reversed, preserving information throughout the process. To simulate arbitrary classical circuits in a reversible form, it is sufficient to construct reversible versions of universal gates such as the NAND gate. Once a reversible version of a universal gate is available, the entire classical computation can be lifted into a reversible framework, which can then be efficiently embedded into a quantum circuit using unitary operations.

**Example 5.1** (Toffoli is universal for classical computation). All boolean logic can be implemented using only NAND gates. NAND and FANOUT (i.e., making a copy of a classical bit  $x$ ) are together universal for classical computation. The Toffoli gate is a controlled-controlled-NOT gate, and with an ancilla initialized to  $|0\rangle$  it computes  $x$  AND  $y$  into the target register. We can use the Toffoli gate to simulate NAND and FANOUT. Therefore the Toffoli gate is universal for classical computation.

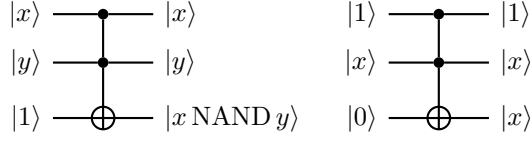


FIGURE 5.1. Using the Toffoli gate to implement NAND and FANOUT

◇

**Exercise 5.1.** Give explicit expressions for using Toffoli gates to implement AND, NOT, XOR, and OR.

A classical computation procedure can be expressed as the evaluation of a boolean map  $f : \{0, 1\}^n \rightarrow \{0, 1\}^m$ , which may be irreversible. However, it can be made into a reversible classical gate

$$(5.1) \quad (z, x) \mapsto (z \oplus f(x), x).$$

In particular,  $(0^m, x) \mapsto (f(x), x)$  is a reversible map that can then be implemented using unitary operations. Efficient implementation of  $x \mapsto f(x)$  on a classical computer means that the number of elementary classical gates (e.g., AND, NOT, NAND gates) is at most  $\text{poly}(n)$ , and the classical implementation of the map uses at most  $\text{poly}(n)$  additional bits for storage. By converting each of the elementary classical gate into a reversible gate, we can implement

$$(5.2) \quad U_f : |0\rangle^{\otimes w} |0\rangle^{\otimes m} |x\rangle \mapsto |g(x)\rangle |f(x)\rangle |x\rangle.$$

Using  $w = \text{poly}(n)$  ancilla qubits, the depth of the quantum circuit is  $\text{poly}(n)$ .

**THEOREM 5.2.** *Any irreversible classical computation using  $\text{poly}(n)$  classical gates can be simulated on a quantum computer using  $\text{poly}(n)$  simple quantum gates and  $\text{poly}(n)$  qubits.*

Up to a polynomial slowdown, a quantum computer is at least as powerful as classical computers. It should be noted that such a procedure is likely to be extremely inefficient. Thus the construction used in Theorem 5.2 is not expected to be practically useful beyond the simplest scenario.

## 5.2. Uncomputation

Unlike classical bits, qubits can exist in superpositions of computational basis states, which enables interference effects in computation. However, qubits are also prone to interference and can easily lose their coherence, causing computational errors. When a quantum computer performs a computation, it can create a large number of ancilla qubits (also called working qubits, or garbage register) that are entangled with the qubits carrying the actual result of the computation. If these ancilla qubits are not properly reset back to their initial state (usually  $|0\rangle^{\otimes a}$ ), they can interfere with subsequent computations and cause errors. This resetting process is called **uncomputation**. Other than avoiding interference, uncomputation is also important for the purpose of resource management. Quantum systems available today have a limited number of qubits. By uncomputing, we can reuse qubits more efficiently.

Uncomputation needs to be done in a very specific way to maintain the integrity of the quantum computation. Simply resetting qubits (for example, by measuring the ancilla qubits and resetting them to  $|0\rangle$ ) is not sufficient, as it can destroy the superposition and entanglement of the other



qubits in the system. Furthermore, due to the no-deleting theorem, there is no generic unitary operator that can set a black-box state to  $|0\rangle^{\otimes w}$ .

Let us now consider how to perform uncomputation when implementing a classical mapping. In quantum computing, an **oracle** means a black box operation that for a given input provides an output, usually the result of evaluating a function on that input. With the help of a working register, we assume that the oracle implementing Eq. (5.2) is available.

In order to set the working register back to  $|0\rangle^{\otimes w}$  while keeping the input and output state, we must use the information stored in  $U_f$  explicitly. We introduce yet another  $m$ -qubit ancilla register initialized at  $|0\rangle^{\otimes m}$ . Then we can use an  $m$ -qubit CNOT controlled on the output register and obtain

$$(5.3) \quad |0\rangle^{\otimes m} |g(x)\rangle |f(x)\rangle |x\rangle \mapsto \underbrace{|f(x)\rangle}_{\text{ancilla}} \underbrace{|g(x)\rangle}_{\text{working}} \underbrace{|f(x)\rangle}_{\text{output}} \underbrace{|x\rangle}_{\text{input}}.$$

It is important to remember that in the operation above, the multi-qubit CNOT gate only performs the classical copying operation in the computational basis, and does not violate the no-cloning theorem.

Recall that  $U_f^{-1} = U_f^\dagger$ , so

$$(5.4) \quad (I^{\otimes m} \otimes U_f^\dagger) |f(x)\rangle |g(x)\rangle |f(x)\rangle |x\rangle = |f(x)\rangle |0\rangle^{\otimes w} |0\rangle^{\otimes m} |x\rangle.$$

Finally we apply an  $m$ -qubit SWAP operator on the ancilla and output registers to obtain

$$(5.5) \quad |f(x)\rangle |0\rangle^{\otimes w} |0\rangle^{\otimes m} |x\rangle \mapsto |0\rangle^{\otimes m} |0\rangle^{\otimes w} |f(x)\rangle |x\rangle.$$

After this procedure, both the ancilla and the working register are set to the initial state. They are no longer entangled to the input or output register, and can be reused for other purposes. The circuit for this uncomputation step is shown in Fig. 5.2.

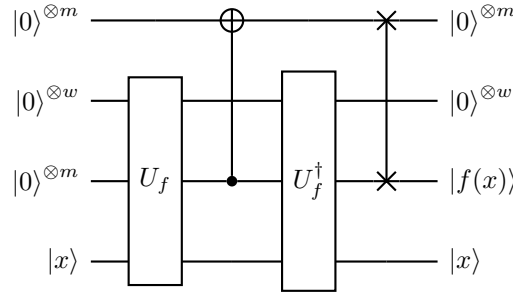


FIGURE 5.2. Circuit for uncomputation. The CNOT and SWAP operators indicate the multi-qubit copy and swap operations, respectively.

**Remark 5.3** (Discarding working registers). After the uncomputation as shown in Fig. 5.2, the first two registers are unchanged before and after the application of the circuit (though they are changed during the intermediate steps). Therefore Fig. 5.2 effectively implements a unitary

$$(5.6) \quad (I^{\otimes(m+w)} \otimes V_f) |0\rangle^{\otimes m} |0\rangle^{\otimes w} |0\rangle^{\otimes m} |x\rangle = |0\rangle^{\otimes m} |0\rangle^{\otimes w} |f(x)\rangle |x\rangle,$$

or equivalently

$$(5.7) \quad V_f |0\rangle^{\otimes m} |x\rangle = |f(x)\rangle |x\rangle.$$

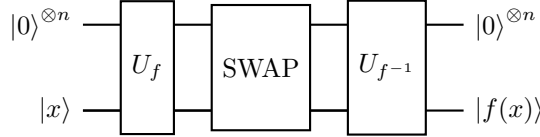
In the definition of  $V_f$ , all working registers have been discarded. This allows us to simplify the notation and focus on the essence of the quantum algorithms under study. Using the technique of uncomputation, if the map  $x \mapsto f(x)$  can be efficiently implemented on a classical computer, then we can implement this map efficiently on a quantum computer as well with a controllable amount of quantum resources.  $\diamond$

**Example 5.4.** Given  $f : \{0, 1\}^n \rightarrow \{0, 1\}^n$ , in general, the transformation  $|x\rangle \mapsto |f(x)\rangle$  is not unitary. However, when  $f$  is a bijection, and we have access to both  $f, f^{-1}$  as follows:

$$(5.8) \quad U_f : |z\rangle |x\rangle \mapsto |z \oplus f(x)\rangle |x\rangle, \quad U_{f^{-1}} : |z\rangle |x\rangle \mapsto |z \oplus f^{-1}(x)\rangle |x\rangle,$$

we can use them to construct the unitary transformation  $U'_f : |x\rangle \mapsto |f(x)\rangle$ .

To implement  $U'_f$ , we will use an ancilla register initialized in the  $|0\rangle^{\otimes n}$  state to hold the result of applying  $f$  or  $f^{-1}$ . Apply  $U_f$  to the state  $|0\rangle^{\otimes n} |x\rangle$  to get  $|f(x)\rangle |x\rangle$ . This setup now contains the desired mapping in the first register, but it is entangled with the input in the second register. Next apply SWAP to the two registers and the state becomes  $|x\rangle |f(x)\rangle$ . Apply  $U_{f^{-1}}$  to the state  $|x\rangle |f(x)\rangle$  to get  $|x \oplus f^{-1}(f(x))\rangle |f(x)\rangle = |x \oplus x\rangle |f(x)\rangle = |0\rangle^{\otimes n} |f(x)\rangle$ . The ancilla register is restored to  $|0\rangle^{\otimes n}$  and can be discarded. This gives our desired  $U'_f$ . The circuit is as follows.



$\diamond$

**Example 5.5.** Another common usage of the uncomputation is to disentangle two registers. Consider the following sequence of operations

$$(5.9) \quad \sum_j c_j |v_j\rangle |0\rangle^{\otimes a} |0\rangle^{\otimes b} \xrightarrow{U_a} \sum_j c_j |v_j\rangle |u_j\rangle |0\rangle^{\otimes b} \xrightarrow{U_b} \sum_j c_j |v_j\rangle |u_j\rangle (\beta_j |0\rangle^{\otimes b} + \sqrt{1 - |\beta_j|^2} |\perp_j\rangle).$$

Here  $U_a$  only acts on the first and second register,  $U_b$  only acts on the second and third register, and  $|\perp_j\rangle$  is a state that is orthogonal to  $|0\rangle^{\otimes b}$ . Our goal is to obtain a state proportional to

$$(5.10) \quad \sum_j c_j \beta_j |v_j\rangle |0\rangle^{\otimes a} |0\rangle^{\otimes b}.$$

This cannot be done by measuring the third register and check whether the outcome is  $0^b$ , since it will lead to  $\sum_j c_j \beta_j |v_j\rangle |u_j\rangle |0\rangle^{\otimes b}$ , which entangles the first two registers. The correct procedure is to perform uncomputation by applying  $U_a^\dagger$  to the first two registers, which gives a state

$$(5.11) \quad \sum_j c_j |v_j\rangle |0\rangle^{\otimes a} (\beta_j |0\rangle^{\otimes b} + \sqrt{1 - |\beta_j|^2} |\perp_j\rangle).$$

Then measuring the third register produces the desired state.  $\diamond$

### 5.3. Fixed point number representation and quantum random access memory

When we want to perform arithmetic operations on a quantum computer, such as addition, multiplication, or more complex functions, we need to encode the numbers we are working with into qubit states. On classical computers, floating point number representations are an efficient way to represent numbers with a wide numerical range. However, on quantum computers, it is often convenient to encode numbers into amplitudes or phases (e.g., via phase kickback). Therefore it is difficult in general to handle numbers that are too large or too small (e.g.,  $3.14 \times 10^{\pm 12}$ ). The standard practice is to use a binary fixed point representation of real numbers.

Any integer  $k \in [N]$  where  $N = 2^n$  can be expressed as an  $n$ -bit string as  $k = (k_{n-1} \cdots k_0)$  with  $k_i \in \{0, 1\}$ . This is called the binary representation of the integer  $k$ . It should be interpreted as

$$(5.12) \quad k = \sum_{i=0}^{n-1} k_i 2^i.$$

The number  $k$  divided by  $2^m$  ( $0 \leq m \leq n$ ) can be written as (note that the binary point is shifted to be after  $k_m$ ):

$$(5.13) \quad a = \frac{k}{2^m} = \sum_{i=0}^{n-1} k_i 2^{i-m} =: (k_{n-1} \cdots k_m . k_{m-1} \cdots k_0).$$

The most common case is  $m = n$ , where

$$(5.14) \quad a = \frac{k}{2^n} = \sum_{i=0}^{n-1} k_i 2^{i-n} =: (0.k_{n-1} \cdots k_0).$$

Sometimes we may also write  $a = 0.k_1 \cdots k_n$ , which is simply a relabeling of the digits. For a given real number  $0 \leq a < 1$  written as

$$(5.15) \quad a = (0.k_1 \cdots k_n k_{n+1} \cdots),$$

the number  $(0.k_1 \cdots k_n)$  is called the  $n$ -bit **fixed point representation** (or  $n$ -bit binary representation) of  $a$ . Therefore to represent  $a$  to additive precision  $\epsilon$ , we will need  $n = \lceil \log_2(1/\epsilon) \rceil$  bits of precision. If the sign of  $a$  is also important, we may reserve one extra bit  $s \in \{0, 1\}$  to indicate its sign and interpret  $(s.k_1 \cdots k_n)$  as  $(-1)^s (0.k_1 \cdots k_n)$ . A complex number  $z$  can be represented using two real numbers as  $z = a + ib$ , where  $a, b \in \mathbb{R}$  are given in the fixed point number representation.

**Definition 5.6.** For a length  $N = 2^n$  classical data vector  $x$ , assume that each component  $x_i$  has a  $d$ -bit representation. Then the **quantum random access memory (QRAM)** is a unitary  $U_{\text{QRAM}}$  acting on  $n + d$  qubits:

$$(5.16) \quad U_{\text{QRAM}} |i\rangle |y\rangle = |i\rangle |y \oplus x_i\rangle.$$

The implementation of  $U_{\text{QRAM}}$  often uses working registers, and such a dependence is hidden in Eq. (5.16) after the uncomputation step. Sometimes QRAM is called the quantum random access classical memory (QRACM). Ideally, the cost for implementing QRAM is  $\text{poly}(n)$ , but this may not be possible if  $x$  represents an unstructured classical data set, and the cost for implementing QRAM may be as high as  $\text{poly}(N)$ .

#### 5.4. Classical arithmetic operations

Using the fixed point number representation and reversible computation, we can approximately implement classical arithmetic operations on quantum computers. The map  $x \mapsto f(x)$  can be implemented as  $U_f |\tilde{x}\rangle |y\rangle = |\tilde{x}\rangle |y \oplus \tilde{f}(x)\rangle$  using e.g., a QRAM. Here  $\tilde{x}$  and  $\tilde{f}(x)$  are  $n$ -bit fixed point representation of  $x, f(x)$  in the computational basis of the quantum register, respectively. However, it may be much more efficient to implement certain classical arithmetic operations on-the-fly on quantum computers without referring to a QRAM. For instance,  $x \mapsto 2x$  can be implemented as a shift operation in the binary format that can be implemented via a sequence of SWAP gates. Other arithmetic mappings, such as  $x \mapsto x^2$ , as well as binary operations  $(x, y) \mapsto x + y$ ,  $(x, y) \mapsto xy$  are harder to implement. Furthermore, these operations can be implemented on quantum computers without going through the process of the reversible implementation of elementary classical gates. Some other classical functions, such as  $x \mapsto \arccos(x)$  can be even more difficult to implement. In general, implementation of classical arithmetic operations on quantum computers will incur a significant overhead, both in terms of the number of ancilla qubits and the circuit depth.

Many arithmetic operations involve a procedure called the controlled rotation, which transforms the information stored in a register from a fixed point representation to the amplitude of the wavefunction.

**Proposition 5.7** (Controlled rotation given rotation angles). *Let  $0 \leq \theta < 1$  have exact  $d$ -bit fixed point representation  $\theta = (. \theta_{d-1} \cdots \theta_0)$ . Then there is a  $(d+1)$ -qubit unitary  $U_\theta$  such that*

$$(5.17) \quad U_\theta : |0\rangle |\theta\rangle \mapsto (\cos(\pi\theta)|0\rangle + \sin(\pi\theta)|1\rangle) |\theta\rangle.$$

PROOF. First (by e.g. Taylor expansion)

$$(5.18) \quad \exp(-i\tau\sigma_y) = \begin{pmatrix} \cos(\tau) & -\sin(\tau) \\ \sin(\tau) & \cos(\tau) \end{pmatrix} =: R_y(2\tau).$$

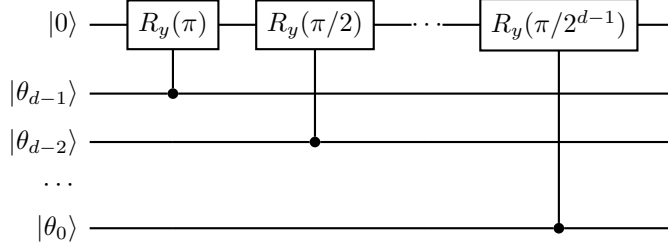
Here  $R_y(\cdot)$  performs a single-qubit rotation around the  $y$ -axis. For any  $j \in [2^d]$  with its binary representation  $j = j_{d-1} \cdots j_0$ , we have

$$(5.19) \quad j/2^d = (.j_{d-1} \cdots j_0).$$

So choose  $\tau = \pi(.j_{d-1} \cdots j_0)$ , and define

$$(5.20) \quad U_\theta = \sum_{j \in [2^d]} \exp(-i\pi(.j_{d-1} \cdots j_0)\sigma_y) \otimes |j\rangle\langle j|.$$

Applying  $U_\theta$  to  $|0\rangle |\theta\rangle$  gives the desired results. This is a sequence of single-qubit rotations on the signal qubit, each controlled by a single qubit.  $\square$

FIGURE 5.3. Quantum circuit for the controlled rotation operation  $U_\theta$ .

**Example 5.8** (Diagonal matrix multiplication using controlled rotation). Let  $0 \leq a < 1$  be given by an  $d$ -bit fixed point representation using an  $d$ -qubit register,  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a function satisfying  $|f(a)| \leq 1$  for all  $0 \leq a < 1$ . For simplicity assume  $f(a) \geq 0$ ; the case of signed  $f(a)$  can be handled by additionally computing the sign of  $f(a)$  and applying a controlled phase flip on the  $|1\rangle$  branch. We would like to construct a circuit that approximately implements

$$(5.21) \quad |a\rangle \rightarrow f(a) |a\rangle.$$

More generally, the state  $|\psi\rangle = \sum_a \psi_a |a\rangle$  is mapped to  $\sum_a \psi_a f(a) |a\rangle$ . This can be viewed as multiplying a diagonal matrix  $D = \text{diag}\{f(a)\}$  to  $|\psi\rangle$ .

To implement such a mapping, we first define

$$(5.22) \quad \theta(a) = \frac{1}{\pi} \arcsin f(a).$$

Note that even though  $a$  is exactly given by  $d$ -bits,  $\theta(a)$  may not be. So we assume that it can be *rounded* to an  $d'$ -bit number  $\tilde{\theta}(a)$ . For simplicity we assume  $d'$  is large enough so that the error of the fixed point representation is negligible in this step. To implement the mapping  $a \mapsto \tilde{\theta}(a)$ , we can construct a classical arithmetics circuit

$$(5.23) \quad U_{\text{angle}} |0^{d'-d}\rangle |a\rangle = |\tilde{\theta}(a)\rangle,$$

whose construction may require  $\text{poly}(m)$  gates and an additional working register of  $\text{poly}(m)$  qubits, which are not displayed here. Therefore the entire controlled rotation operation needed is given by the circuit in Fig. 5.4.

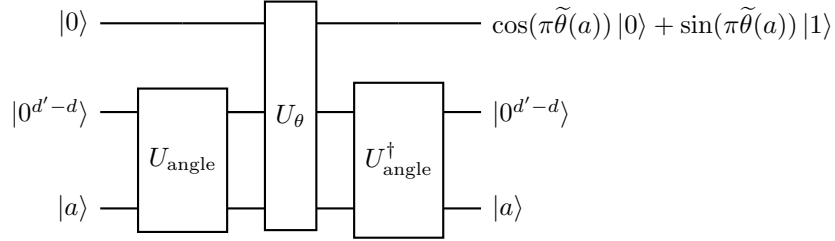


FIGURE 5.4. Circuit for using controlled rotation to implement the multiplication of a diagonal matrix (not including additional working register for classical arithmetic operations).

Note that through the uncomputation  $U_{\text{angle}}^\dagger$ , the  $d' - d$  ancilla qubits also become a working register. After uncomputation, the ancillas are returned to  $|0^{d'-d}\rangle$  (together with any additional workspace used in  $U_{\text{angle}}$ ), so they may be reused. We obtain a unitary  $U_{\text{CR}}$  satisfying

$$(5.24) \quad U_{\text{CR}} |0\rangle |a\rangle = \left( \cos(\pi\tilde{\theta}(a)) |0\rangle + \sin(\pi\tilde{\theta}(a)) |1\rangle \right) |a\rangle \approx \left( \sqrt{1 - f(a)^2} |0\rangle + f(a) |1\rangle \right) |a\rangle.$$

Measure the single ancilla qubit. If the result is 1, the data register is projected onto a state proportional to  $\sum_a \psi_a f(a) |a\rangle$ , i.e., the mapping in Eq. (5.21) up to renormalization. If the input state is  $|\psi\rangle = \sum_a \psi_a |a\rangle$ , the probability of obtaining 1 after measuring the ancilla qubit is

$$(5.25) \quad \mathbb{P}(1) \approx \sum_a |\psi_a|^2 |f(a)|^2.$$

◇

**Example 5.9** (Use of arithmetic operations in the HHL algorithm). The last step of the Harrow–Hassidim–Lloyd (HHL) algorithm for solving a linear system of equations  $Ax = b$  with a Hermitian matrix  $A$  involves the following arithmetic operations. For simplicity assume  $\lambda_j$  (eigenvalues of  $A$ ) are given exactly in a  $d$ -bit fixed point number representation, and  $\lambda_j \in [\delta, 1]$  for some  $\delta > 0$ . Start from a linear combination of states  $|\psi\rangle = \sum_j \beta_j |0\rangle |\lambda_j\rangle |v_j\rangle$ , we would like to construct a state

$$(5.26) \quad |\psi'\rangle = \sum_j \frac{C\beta_j}{\lambda_j} |1\rangle |\lambda_j\rangle |v_j\rangle + |0\rangle |\perp\rangle.$$

Here  $C$  is a normalization constant chosen so that  $|C/\lambda_j| < 1$  for all  $\lambda_j \in [\delta, 1]$ , and  $|\perp\rangle$  is an irrelevant unnormalized state. Viewing this as a diagonal matrix multiplication problem, the function of interest is

$$(5.27) \quad f(a) = \frac{C}{a}, \quad a \in [\delta, 1].$$

The implementation involves the classical arithmetic circuit for computing

$$(5.28) \quad \theta(a) = \frac{1}{\pi} \arcsin f(a) = \frac{1}{\pi} \arcsin(C/a)$$

using  $d'$  bits ( $d' > d$ ).

Once  $|\psi'\rangle$  is prepared, we can uncompute  $|\lambda_j\rangle$  to obtain a state

$$(5.29) \quad \sum_j \frac{C\beta_j}{\lambda_j} |1\rangle |0^d\rangle |v_j\rangle + |0\rangle |\perp'\rangle$$

to disentangle the  $\lambda_j$  register from the  $v_j$  register. If we measure the first ancilla register and obtain 1, we obtain the desired form of the solution in the HHL algorithm.  $\diamond$

### Notes and further reading

Reversible computation predates quantum computing and has both physical and algorithmic motivations. Landauer related logical irreversibility to dissipation [Lan61]. For background on reversible embeddings of classical circuits into unitary dynamics, see [NC00, Section 3.2.5]. For fixed-point encodings and reversible arithmetic (addition, multiplication, and function evaluation), a detailed treatment is given in [RP11, Chapter 6]. For standard universal classical gate constructions and decompositions into elementary quantum gates, see [BBC<sup>+</sup>95]. There is also opportunity to optimize the cost of the uncomputation stage. An example is Gidney's construction [Gid18] of the quantum adder circuit. The QRAM model [GLM08] should be interpreted as an assumption about data access rather than an automatic feature of a fault-tolerant architecture.





## CHAPTER 6

# Query complexity and quantum complexity theory



## CHAPTER 7

# **Perturbation theory**



## CHAPTER 8

# Statistical estimates



Part III

Algorithm





## CHAPTER 9

# Block encoding

This chapter introduces block encoding as an input model for matrix problems on a quantum computer. The basic difficulty is that many tasks in scientific computation are naturally phrased in terms of non-unitary linear maps, whereas the native operations available to quantum hardware are unitary. Block encoding addresses this mismatch by representing a target matrix  $A$  (up to a subnormalization factor and a prescribed error tolerance) as a distinguished submatrix block of a larger unitary  $U_A$ , so that applying  $U_A$  and post-selecting on ancilla qubits effectively applies  $A$  to a state.

The possibility of constructing an efficiently implementable  $U_A$  depends strongly on the structure of  $A$  and on the assumed access model. For a dense matrix without additional structure, any reasonable input model is typically prohibitive, since the input description may itself require exponential resources. We therefore focus on representative settings in which block encodings can be constructed efficiently under suitable oracle access assumptions.

The true power of block encoding does not come directly from the ability to represent arbitrary matrices within blocks of a larger unitary. Rather, it stems from the ability to compose block encodings to block encode more complicated matrices and functions of matrices. We then describe how block encodings can be combined to obtain encodings of matrix additions and multiplications, while tracking the corresponding subnormalization factors and errors. Linear combinations of unitaries provide a flexible mechanism for such constructions. In this way, block encoding serves as an interface between matrix-oriented problem statements and unitary circuit realizations used throughout subsequent chapters.

### 9.1. Block encoding

The simplest example of block encoding is the following: assume we can find a  $(n + 1)$ -qubit unitary matrix  $U_A \in \text{U}(2N)$  (where  $N = 2^n$ ) such that

$$U_A = \begin{pmatrix} A & * \\ * & * \end{pmatrix}$$

where  $*$  means that the corresponding matrix entries are irrelevant, then for any  $n$ -qubit quantum state  $|b\rangle$ , we can consider the state

$$(9.1) \quad |0, b\rangle = |0\rangle |b\rangle = \begin{pmatrix} b \\ 0 \end{pmatrix},$$

and

$$(9.2) \quad U_A |0, b\rangle = \begin{pmatrix} Ab \\ * \end{pmatrix} =: |0\rangle A|b\rangle + |\perp\rangle.$$

Here the (unnormalized) state  $|\perp\rangle$  can be written as  $|1\rangle|\psi\rangle$  for some (unnormalized) state  $|\psi\rangle$  that is irrelevant to the computation of  $A|b\rangle$ . In particular, it satisfies the orthogonality relation.

$$(9.3) \quad (|0\rangle\langle 0| \otimes I_N) |\perp\rangle = 0.$$

In order to obtain  $A|b\rangle$ , we measure the ancilla qubit and postselect on the outcome 0. This can be summarized into the following quantum circuit:

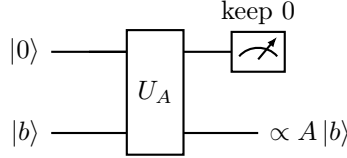


FIGURE 9.1. Circuit for block encoding of  $A$  using one ancilla qubit. By measuring the ancilla qubit and postselecting on the outcome 0, the state in the system register is a normalized state proportional to  $A|b\rangle$ .

Note that the output state is normalized after the measurement takes place. The success probability of obtaining 0 from the measurement can be computed as

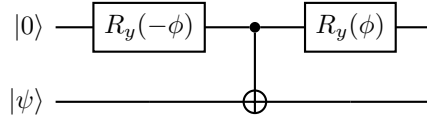
$$(9.4) \quad p(0) = \|A|b\rangle\|^2 = \langle b|A^\dagger A|b\rangle.$$

So the missing information of the norm  $\|A|b\rangle\|$  can be recovered via the success probability  $p(0)$  if needed. We find that the success probability is only determined by  $A, |b\rangle$ , and is independent of other irrelevant components of  $U_A$ .

**Example 9.1.** Consider the  $2 \times 2$  matrix

$$(9.5) \quad A = \frac{3}{4}I + \frac{1}{4}X = \begin{pmatrix} 0.75 & 0.25 \\ 0.25 & 0.75 \end{pmatrix}.$$

Consider the following circuit ( $\phi = \frac{1}{3}\pi$ )



Here

$$(9.6) \quad R_y(\theta) := \begin{bmatrix} \cos\left(\frac{\theta}{2}\right) & -\sin\left(\frac{\theta}{2}\right) \\ \sin\left(\frac{\theta}{2}\right) & \cos\left(\frac{\theta}{2}\right) \end{bmatrix} = e^{-i\theta Y/2}$$

is the  $Y$ -rotation matrix. One may directly verify that  $U_A$  is an exact block encoding of  $A$  using one ancilla qubit.  $\diamond$

Note that we may not need to restrict the matrix  $U_A$  to be an  $(n+1)$ -qubit matrix. If we can find any  $(n+m)$ -qubit unitary matrix  $U_A$  so that

$$(9.7) \quad U_A = \begin{pmatrix} A & * & \cdots & * \\ * & * & \cdots & * \\ \vdots & & \ddots & \\ * & * & \cdots & * \end{pmatrix}$$

Here each  $*$  stands for an  $n$ -qubit matrix, and there are  $2^m$  block rows / columns in  $U_A$ . Using the partial application of operators in Definition 2.25, the relation above can be written compactly using the bracket notation as

$$(9.8) \quad A = \langle 0^m | U_A | 0^m \rangle.$$

**Exercise 9.1.** Given a unitary matrix  $U$  and any submatrix block  $A$ , prove that  $\|A\| \leq 1$ .

In order to find such a block encoding  $U_A$ , Exercise 9.1 shows that a necessary condition for the existence of  $U_A$  is that  $\|A\| \leq 1$ . However, if we can find sufficiently large  $\alpha$  and  $U_A$  so that

$$(9.9) \quad A/\alpha = \langle 0^m | U_A | 0^m \rangle.$$

By measuring the  $m$  ancilla qubits and postselecting on the outcome  $0^m$ , we still obtain the normalized state  $\frac{A|b\rangle}{\|A|b\rangle\|}$ . The number  $\alpha$  is hidden in the success probability:

$$(9.10) \quad p(0^m) = \frac{1}{\alpha^2} \|A|b\rangle\|^2 = \frac{1}{\alpha^2} \langle b | A^\dagger A | b \rangle.$$

So if  $\alpha$  is chosen to be too large, the probability of obtaining all 0's from the measurement can be vanishingly small.

Finally, it can be difficult to find  $U_A$  to block encode  $A$  exactly. This is not a problem, since it is sufficient if we can find  $U_A$  to block encode  $A$  up to some error  $\epsilon$ . We are now ready to give the definition of block encoding in Definition 9.2.

**Definition 9.2 (Block encoding).** Given an  $n$ -qubit matrix  $A$ , if we can find  $\alpha, \epsilon \in \mathbb{R}_+$ , and an  $(m+n)$ -qubit unitary matrix  $U_A$  so that

$$(9.11) \quad \|A - \alpha \langle 0^m | U_A | 0^m \rangle\| \leq \epsilon,$$

then  $U_A$  is called an  $(\alpha, m, \epsilon)$ -block-encoding of  $A$ . When the block encoding is exact with  $\epsilon = 0$ ,  $U_A$  is called an  $(\alpha, m)$ -block-encoding of  $A$ . The set of all  $(\alpha, m, \epsilon)$ -block-encodings of  $A$  is denoted by  $\text{BE}_{\alpha, m}(A, \epsilon)$ . The parameter  $\alpha$  is referred to as the block encoding factor, or the subnormalization factor.

When discussing block encodings, we often ignore certain errors such as the error in the finite precision number representation. We define a shorthand notation  $\text{BE}_{\alpha, m}(A) = \text{BE}_{\alpha, m}(A, 0)$ . Assume we know each matrix element of the  $n$ -qubit matrix  $A_{ij}$ , and we are given an  $(n+m)$ -qubit unitary  $U_A$ . In order to verify that  $U_A \in \text{BE}_{1, m}(A)$ , we only need to verify that

$$(9.12) \quad \langle 0^m, i | U_A | 0^m, j \rangle = A_{ij},$$

and  $U_A$  applied to any vector  $|0^m, b\rangle$  can be obtained via the superposition principle.

Therefore we may first evaluate the state  $U_A |0^m, j\rangle$ , perform an inner product with  $|0^m, i\rangle$ , and verify the resulting inner product is  $A_{ij}$ . We will also use the following technique frequently. Assume  $U_A = U_B U_C$ , and then

$$(9.13) \quad \langle 0^m, i | U_A | 0^m, j \rangle = \langle 0^m, i | U_B U_C | 0^m, j \rangle = (U_B^\dagger | 0^m, i \rangle)^\dagger (U_C | 0^m, j \rangle).$$

So we can evaluate the states  $U_B^\dagger |0^m, i\rangle, U_C |0^m, j\rangle$  independently, and then verify the inner product is  $A_{ij}$ . Such a calculation amounts to running the circuit Fig. 9.2, and if the ancilla qubits are measured to be  $0^m$ , the system qubits return the normalized state  $\sum_i A_{ij} |i\rangle / \|\sum_i A_{ij} |i\rangle\|$ .

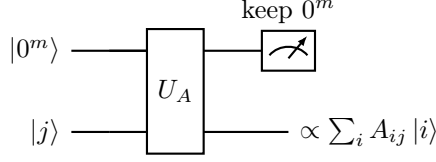


FIGURE 9.2. Circuit for general block encoding of  $A$ .

**Example 9.3.** For any  $n$ -qubit matrix  $A$  with  $\|A\| \leq 1$  with singular value decomposition  $A = W\Sigma V^\dagger$ , all singular values in the diagonal matrix  $\Sigma$  are in  $[0, 1]$ . Then we may construct an  $(n+1)$ -qubit unitary matrix ( $N = 2^n$ )

$$(9.14) \quad \begin{aligned} U_A &:= \begin{pmatrix} W & 0 \\ 0 & I_N \end{pmatrix} \begin{pmatrix} \Sigma & \sqrt{I_N - \Sigma^2} \\ \sqrt{I_N - \Sigma^2} & -\Sigma \end{pmatrix} \begin{pmatrix} V^\dagger & 0 \\ 0 & I_N \end{pmatrix} \\ &= \begin{pmatrix} A & W\sqrt{I_N - \Sigma^2} \\ \sqrt{I_N - \Sigma^2}V^\dagger & -\Sigma \end{pmatrix} \end{aligned}$$

which is a  $(1, 1)$ -block-encoding of  $A$ .  $\diamond$

Example 9.3 shows that in principle, any matrix  $A$  with  $\|A\| \leq 1$  can be accessed via a  $(1, 1, 0)$ -block-encoding. However, this construction does not state how to construct  $A$  using simple one and two qubit gates.

**Example 9.4** (Random circuit block encoded matrix). How can we construct a pseudo-random non-unitary matrix on a quantum computer? A naive approach would be to generate a dense pseudo-random matrix  $A$  classically and then encode it into a quantum circuit. However, this is highly inefficient in practice, particularly for large matrices, due to the exponential overhead in loading dense classical data into a quantum system.

Instead, we seek to work with matrices that are inherently easy to generate within a quantum circuit model. This motivates the **random circuit based block-encoded matrix** (RACBEM) model. Rather than first constructing a matrix  $A$  and then searching for a block-encoding unitary  $U_A$ , the RACBEM model reverses the thought process: we begin by constructing a unitary  $U_A$  that is easy to implement on a quantum computer, typically using random quantum circuits, and then extract  $A$  as a subblock of  $U_A$ . This provides a practical and scalable way to generate structured pseudo-random non-unitary matrices compatible with quantum algorithm design. Similar to the LINPACK benchmark, which is used to rank classical supercomputers in the TOP500 list by solving  $Ax = b$  for pseudorandom matrices  $A$ , such block-encoded pseudorandom matrices can serve as a useful tool for benchmarking scientific computing applications on quantum computers.  $\diamond$

## 9.2. Linear combination of unitaries

The **linear combination of unitaries** (LCU) is an important quantum primitive, which allows quantum algorithms to be implemented as a superposition of unitary matrices rather than

attempting to find a single unitary that accomplishes a desired task. This often simplifies the design and analysis of quantum algorithms. LCU can also be viewed as a special way for constructing block encoding. Combined with a technique called qubitization, which will be discussed in detail in ??, LCU can be used to implement a large class of matrix functions (eigenvalue transformations) and generalized matrix functions (singular value transformations).

Let  $T = \sum_{i=0}^{K-1} \alpha_i U_i$  be a linear combination of unitary matrices  $U_i$ . For simplicity let  $K = 2^a$ . Then

$$(9.15) \quad U_{\text{SEL}} := \sum_{i \in [K]} |i\rangle \langle i| \otimes U_i,$$

implements the selection of  $U_i$  conditioned on the value of the  $a$ -qubit ancilla register (also called the control register).  $U_{\text{SEL}}$  is called a **select oracle**.

If all linear combination coefficients  $\alpha_i \geq 0$ , we can let  $V_{\text{PREP}}$  be a unitary operation satisfying

$$(9.16) \quad V_{\text{PREP}} |0^a\rangle = \frac{1}{\sqrt{\|\alpha\|_1}} \sum_{i \in [K]} \sqrt{\alpha_i} |i\rangle,$$

which is called a **prepare oracle**. The 1-norm of the coefficients is given by

$$(9.17) \quad \|\alpha\|_1 = \sum_i |\alpha_i|.$$

In matrix form,

$$(9.18) \quad V_{\text{PREP}} = \frac{1}{\sqrt{\|\alpha\|_1}} \begin{pmatrix} \sqrt{\alpha_0} & * & \cdots & * \\ \vdots & * & \ddots & \vdots \\ \sqrt{\alpha_{K-1}} & * & \cdots & * \end{pmatrix}.$$

where the first column is  $V_{\text{PREP}} |0^a\rangle$ , and all other columns are orthogonal to it. Then

$$(9.19) \quad V_{\text{PREP}}^\dagger = \frac{1}{\sqrt{\|\alpha\|_1}} \begin{pmatrix} \sqrt{\alpha_0} & \cdots & \sqrt{\alpha_{K-1}} \\ * & \cdots & * \\ \vdots & \ddots & \vdots \\ * & \cdots & * \end{pmatrix}.$$

More generally, we can arbitrarily decompose  $\alpha_i = \beta_i \gamma_i$ , so that

$$(9.20) \quad V_{\text{PREP}} = \frac{1}{\|\beta\|_2} \begin{pmatrix} \beta_0 & * & \cdots & * \\ \vdots & * & \ddots & \vdots \\ \beta_{K-1} & * & \cdots & * \end{pmatrix}, \quad \tilde{V}_{\text{PREP}} = \frac{1}{\|\gamma\|_2} \begin{pmatrix} \gamma_0 & \cdots & \gamma_{K-1} \\ * & \cdots & * \\ \vdots & \ddots & \vdots \\ * & \cdots & * \end{pmatrix}$$

are unitaries and can be efficiently implemented. When  $\alpha_i \geq 0$ , we can choose  $\beta_i = \gamma_i = \sqrt{\alpha_i}$  which gives  $\tilde{V}_{\text{PREP}} = V_{\text{PREP}}^\dagger$ . Then  $T$  can be implemented using the unitary given in Lemma 9.5.

**Lemma 9.5** (Linear combination of unitaries). *For*

$$(9.21) \quad T = \sum_{i=0}^{K-1} \alpha_i U_i, \quad \alpha_i = \beta_i \gamma_i, \quad K = 2^a, \quad U_i \in \mathcal{U}(2^n),$$

let  $U_{\text{SEL}}, V_{\text{PREP}}, \tilde{V}_{\text{PREP}}$  be given in Eqs. (9.15) and (9.20), respectively. Define

$$(9.22) \quad W = (\tilde{V}_{\text{PREP}} \otimes I_n) U_{\text{SEL}} (V_{\text{PREP}} \otimes I_n)$$

as implemented in Fig. 9.3. Then  $W \in \text{BE}_{\|\beta\|_2 \|\gamma\|_2, a}(T)$ . The smallest subnormalization factor is obtained by setting

$$(9.23) \quad |\beta_i| = |\gamma_i| = \sqrt{|\alpha_i|}, \quad i \in [K],$$

and  $W \in \text{BE}_{\|\alpha\|_1, a}(T)$ .

PROOF. For any  $n$ -qubit state  $|\psi\rangle$ ,

$$(9.24) \quad U_{\text{SEL}}(V_{\text{PREP}} \otimes I_n) |0^a\rangle |\psi\rangle = U_{\text{SEL}} \frac{1}{\|\beta\|_2} \sum_i \beta_i |i\rangle |\psi\rangle = \frac{1}{\|\beta\|_2} \sum_i \beta_i |i\rangle U_i |\psi\rangle.$$

Let the state  $|\tilde{\perp}\rangle$  collect all the states marked by  $*$  orthogonal to  $|0^a\rangle$ , and use  $\beta_i \gamma_i = \alpha_i$ ,

$$(9.25) \quad (\tilde{V}_{\text{PREP}} \otimes I_n) U_{\text{SEL}}(V_{\text{PREP}} \otimes I_n) |0^a\rangle |\psi\rangle = \frac{1}{\|\beta\|_2 \|\gamma\|_2} |0^a\rangle \sum_i \alpha_i U_i |\psi\rangle + |\tilde{\perp}\rangle = \frac{1}{\|\beta\|_2 \|\gamma\|_2} |0^a\rangle T |\psi\rangle + |\tilde{\perp}\rangle.$$

Use Cauchy-Schwarz

$$(9.26) \quad \|\alpha\|_1 = \sum_i |\alpha_i| = \sum_i |\beta_i \gamma_i| \leq \|\beta\|_2 \|\gamma\|_2,$$

we find that the optimal prepare oracle should satisfy  $|\beta_i| = |\gamma_i| = \sqrt{|\alpha_i|}, \forall i$ .  $\square$

The LCU Lemma states that the number of ancilla qubits needed only depends logarithmically on  $K$ , the number of terms in the linear combination. Hence it is possible to implement the linear combination of a very large number of terms efficiently. From a practical perspective, the select and prepare oracles use multi-qubit controls, and may be difficult to implement. If implemented directly, the number of multi-qubit controls again depends linearly on  $K$  and is not desirable. Therefore an efficient implementation using LCU (in terms of the gate complexity) also requires additional structure in the prepare and select oracles.

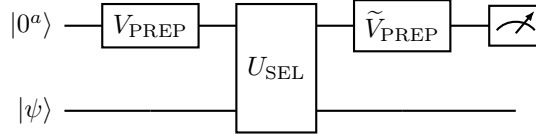


FIGURE 9.3. Circuit for linear combination of unitaries. When all coefficients are nonnegative, we may set  $\tilde{V}_{\text{PREP}} = V_{\text{PREP}}^\dagger$ .

**Example 9.6.** If we apply  $W$  to  $|0^a\rangle |\psi\rangle$  and measure the ancilla qubits, then the probability of obtaining the outcome  $0^a$  in the ancilla qubits (and therefore obtaining the state  $T |\psi\rangle / \|T |\psi\rangle\|$  in the system register) is  $(\|T |\psi\rangle\| / \|\alpha\|_1)^2$ . The expected number of repetition needed to succeed is  $(\|\alpha\|_1 / \|T |\psi\rangle\|)^2$ . Using amplitude amplification (AA) in ??, this number can be reduced to  $\mathcal{O}(\|\alpha\|_1 / \|T |\psi\rangle\|)$ .  $\diamond$

An important application of LCU is that if  $A, B$  can be accessed via their block encodings, then we can construct a block encoding of the matrix addition  $A + B$ .

**Example 9.7** (Linear combination of two block encoded matrices). Let  $U_A, U_B$  be two  $n$ -qubit unitaries, and we would like to construct a block encoding of  $T = U_A + U_B$ .

There are two terms in total, so one ancilla qubit is needed. The prepare oracle needs to implement

$$(9.27) \quad V_{\text{PREP}} |0\rangle = \frac{1}{\sqrt{2}}(|0\rangle + |1\rangle),$$

so this is the Hadamard gate. The circuit is given by Fig. 9.4, which constructs  $W \in \text{BE}_{2,1}(T)$ .

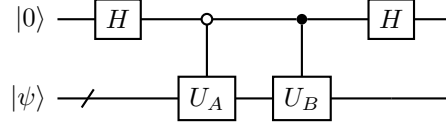


FIGURE 9.4. Circuit for linear combination of two unitaries.

◇

**Exercise 9.2.** Let  $A, B$  be two  $n$ -qubit matrices encoded by  $U_A \in \text{BE}_{1,m}(A), U_B \in \text{BE}_{1,m}(B)$ . Construct a circuit to block encode  $C = A + B$ . What about  $U_A \in \text{BE}_{\alpha_A,m}(A), U_B \in \text{BE}_{\alpha_B,m}(B)$ ?

**Exercise 9.3.** Consider a system described by the linear combination  $T = X + Y + 2Z$ , where  $X, Y, Z$  are the Pauli matrices. Construct a select oracle  $U$  for this system, and describe how to use the LCU technique to construct a block encoding of  $T$ .

**Example 9.8.** Consider the following TFIM model with periodic boundary conditions ( $Z_n = Z_0$ ), and  $n = 2^n$ ,

$$(9.28) \quad \hat{H} = - \sum_{i \in [n]} Z_i Z_{i+1} - \sum_{i \in [n]} X_i.$$

In order to use LCU, we need  $(n+1)$  ancilla qubits. In this case, the prepare oracle can be simply constructed from the Hadamard gate

$$(9.29) \quad V_{\text{PREP}} = H^{\otimes(n+1)},$$

and the select oracle implements

$$(9.30) \quad U_{\text{SEL}} = \sum_{i \in [n]} |i\rangle \langle i| \otimes (-Z_i Z_{i+1}) + \sum_{i \in [n]} |i+n\rangle \langle i+n| \otimes (-X_i).$$

The corresponding  $W \in \text{BE}_{2n, n+1}(\hat{H})$ .

◇

**Example 9.9** (Highly oscillatory integral). Consider evaluating the matrix integral  $\int_0^1 A(s) ds$ , where  $A(s) \in \mathbb{C}^{2^n \times 2^n}$ ,  $A(0) = A(1)$  and  $\sup_{s \in [0,1]} \|A(s)\| \leq 1$ . Given that the entries of  $A(s)$  exhibit significant oscillations as a function of  $s$ , in general there is no known efficient method (classical or quantum) to compute this integral without using a sufficiently fine grid and numerical quadrature. For simplicity, we adopt a uniform grid defined by  $\{s_k = \frac{k}{M}\}_{k=0}^M$ , where  $M$  is sufficiently large, to implement the quadrature method.

$$(9.31) \quad \int_0^1 A(s) ds = \frac{1}{M} \sum_{k=0}^{M-1} A(k/M) + E, \quad \|E\| \leq \epsilon.$$

For each  $s$ , assume that  $A(s)$  has a  $(1, a, 0)$ -block encoding denoted by  $U_{A(s)}$ , and the  $s$ -dependence can be implemented coherently using e.g., classical arithmetic operations. In a discretized setting, let  $M = 2^m$ , this means that the following select oracle defined on a register with  $m + a + n$  qubits:

$$(9.32) \quad U_{\text{SEL}} = \sum_{k=0}^{M-1} |k\rangle\langle k| \otimes U_{A(k/M)},$$

which we assume can be efficiently implemented with cost  $\text{poly}(mn)$ . The prepare oracle is simply the  $m$ -qubit Hadamard gate  $H^{\otimes m}$ . Then the circuit  $(H^{\otimes m} \otimes I_{a+n}) U_{\text{SEL}} (H^{\otimes m} \otimes I_{a+n})$  is a  $(1, a + m, \epsilon)$ -block encoding of the matrix-valued integral  $\int_0^1 A(s) ds$ . It uses  $m$  ancilla qubits, and the gate complexity is dominated by that of the select oracle and is  $\text{poly}(mn)$ . This is an exponential improvement in the parameter  $M$  for constructing such a block encoding, compared to a direct classical quadrature implementation whose cost is at least linear in  $M$ .  $\diamond$

### 9.3. Block encodings of matrix additions and multiplications

We now record basic composition rules for block encodings that will be used throughout the book.

The linear combination of unitaries (LCU) construction from Section 9.2 immediately yields a block encoding of a sum of block-encoded matrices. For simplicity, we state the result for  $M = 2^m$  summands.

**Proposition 9.10** (Sum of  $M$  block-encoded matrices). *Let  $M = 2^m$  and let  $A_0, \dots, A_{M-1}$  be matrices of the same dimension. Assume that for each  $j \in [M]$  we are given a block encoding*

$$(9.33) \quad U_{A_j} \in \text{BE}_{\alpha_j, a}(A_j, \epsilon_j), \quad \alpha_j \geq 0.$$

*Set  $\gamma := \sum_{j=0}^{M-1} \alpha_j > 0$ . Let  $U_{\text{SEL}} := \sum_{j \in [M]} |j\rangle\langle j| \otimes U_{A_j}$  be the select oracle acting on an  $m$ -qubit control register, the  $a$ -qubit ancilla register, and the system register. Let  $V_{\text{PREP}}$  be any unitary on the  $m$ -qubit control register satisfying*

$$(9.34) \quad V_{\text{PREP}} |0^m\rangle = \frac{1}{\sqrt{\gamma}} \sum_{j=0}^{M-1} \sqrt{\alpha_j} |j\rangle.$$

*Define*

$$(9.35) \quad W := (V_{\text{PREP}}^\dagger \otimes I) U_{\text{SEL}} (V_{\text{PREP}} \otimes I).$$

*Then*

$$(9.36) \quad W \in \text{BE}_{\gamma, a+m} \left( \sum_{j=0}^{M-1} A_j, \sum_{j=0}^{M-1} \epsilon_j \right).$$

**PROOF.** Write  $B_j := \langle 0^a | U_{A_j} | 0^a \rangle$ , so that  $\|A_j - \alpha_j B_j\| \leq \epsilon_j$  and  $\|B_j\| \leq 1$ . By direct computation of the  $(|0^m, 0^a\rangle)$  block,

$$(9.37) \quad \langle 0^m, 0^a | W | 0^m, 0^a \rangle = \sum_{j=0}^{M-1} \frac{\alpha_j}{\gamma} B_j.$$



Therefore

$$(9.38) \quad \left\| \sum_{j=0}^{M-1} A_j - \gamma \langle 0^m, 0^a | W | 0^m, 0^a \rangle \right\| \leq \sum_{j=0}^{M-1} \|A_j - \alpha_j B_j\| \leq \sum_{j=0}^{M-1} \epsilon_j,$$

which is the claimed block-encoding statement.  $\square$

We next record a simple (though not always ancilla-optimal) rule for block encoding a product.

**Proposition 9.11** (Product of  $M$  block-encoded matrices). *Let  $A_0, \dots, A_{M-1}$  be matrices with compatible dimensions. Assume that for each  $j \in [M]$  we are given*

$$(9.39) \quad U_{A_j} \in \text{BE}_{\alpha_j, a_j}(A_j, \epsilon_j).$$

*Let  $U$  be the unitary obtained by applying  $U_{A_0}, U_{A_1}, \dots, U_{A_{M-1}}$  sequentially on disjoint ancilla registers (of sizes  $a_0, \dots, a_{M-1}$ ) and a common system register. Then*

$$(9.40) \quad U \in \text{BE}_{\prod_{j=0}^{M-1} \alpha_j, \sum_{j=0}^{M-1} a_j} \left( A_{M-1} \cdots A_0, \prod_{j=0}^{M-1} (\alpha_j + \epsilon_j) - \prod_{j=0}^{M-1} \alpha_j \right).$$

PROOF. For each  $j$ , define  $B_j := \langle 0^{a_j} | U_{A_j} | 0^{a_j} \rangle$  so that  $\|A_j - \alpha_j B_j\| \leq \epsilon_j$  and  $\|B_j\| \leq 1$ . Since the ancilla registers are disjoint, we have

$$(9.41) \quad \langle 0^{a_0+\dots+a_{M-1}} | U | 0^{a_0+\dots+a_{M-1}} \rangle = B_{M-1} \cdots B_0.$$

It remains to bound

$$(9.42) \quad \left\| A_{M-1} \cdots A_0 - \left( \prod_{j=0}^{M-1} \alpha_j \right) B_{M-1} \cdots B_0 \right\|.$$

We prove by induction on  $M$  the inequality

$$(9.43) \quad \left\| \prod_{j=0}^{M-1} A_j - \prod_{j=0}^{M-1} (\alpha_j B_j) \right\| \leq \prod_{j=0}^{M-1} (\alpha_j + \epsilon_j) - \prod_{j=0}^{M-1} \alpha_j.$$

The case  $M = 1$  is immediate. For the induction step, write  $P := \prod_{j=0}^{M-2} A_j$  and  $\tilde{P} := \prod_{j=0}^{M-2} (\alpha_j B_j)$ . Then

$$(9.44) \quad \begin{aligned} \left\| A_{M-1} P - (\alpha_{M-1} B_{M-1}) \tilde{P} \right\| &\leq \left\| (A_{M-1} - \alpha_{M-1} B_{M-1}) P \right\| + \left\| \alpha_{M-1} B_{M-1} (P - \tilde{P}) \right\| \\ &\leq \epsilon_{M-1} \|P\| + \alpha_{M-1} \|P - \tilde{P}\|. \end{aligned}$$

Using  $\|A_j\| \leq \alpha_j + \epsilon_j$  (by  $\|A_j\| \leq \|A_j - \alpha_j B_j\| + \alpha_j \|B_j\|$ ), we have  $\|P\| \leq \prod_{j=0}^{M-2} (\alpha_j + \epsilon_j)$ . Applying the induction hypothesis to  $\|P - \tilde{P}\|$  yields

$$(9.45) \quad \begin{aligned} \left\| A_{M-1} P - (\alpha_{M-1} B_{M-1}) \tilde{P} \right\| &\leq \epsilon_{M-1} \prod_{j=0}^{M-2} (\alpha_j + \epsilon_j) + \alpha_{M-1} \left( \prod_{j=0}^{M-2} (\alpha_j + \epsilon_j) - \prod_{j=0}^{M-2} \alpha_j \right) \\ &= \prod_{j=0}^{M-1} (\alpha_j + \epsilon_j) - \prod_{j=0}^{M-1} \alpha_j, \end{aligned}$$

completing the induction.  $\square$

**Example 9.12** (Multiplication of block encoded matrices). If  $A, B$  are given by their block encodings  $U_A \in \text{BE}_{\alpha,a}(A), U_B \in \text{BE}_{\beta,b}(B)$ , then the product  $AB$  can also be block encoded (see Fig. 9.5), which uses  $a + b$  ancilla qubits. This is because  $AB/(\alpha\beta) = \langle 0^{a+b} | (U_A \otimes I_b)(I_a \otimes U_B) | 0^{a+b} \rangle$ . Hence  $(U_A \otimes I_b)(I_a \otimes U_B) \in \text{BE}_{\alpha\beta, a+b}(AB)$ .

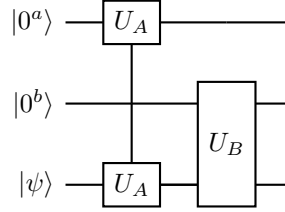


FIGURE 9.5. Quantum circuit for block encoding the product of matrices using  $a + b$  ancilla qubits.

However, this is not the most efficient way for block encoding the product of matrices. In ??, we have demonstrated that using deferred measurement, we only need one extra ancilla qubit to record whether the ancilla register is in the all 0 state. Specifically, assume  $a = b$  for simplicity; Fig. 9.6 is a schematic circuit (the control denotes a check of the ancilla register being in  $|0^a\rangle$ ) that constructs a unitary in  $\text{BE}_{\alpha\beta, a+1}(AB)$ .

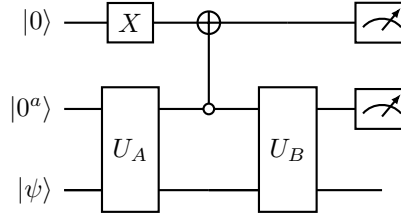


FIGURE 9.6. Quantum circuit for block encoding the product of matrices using  $a + 1$  ancilla qubits (assuming  $a = b$ ).

Following this strategy, when multiplying  $L$  matrices  $A_i$  each given by  $U_{A_i} \in \text{BE}_{\alpha_i, a}(A_i)$ , we can introduce  $L - 1$  ancilla qubits to obtain a unitary in  $\text{BE}_{\prod_{i=1}^L \alpha_i, a+L-1}(A_L \cdots A_1)$ . Even more efficiently, using the compression gadget in ??, the number of ancilla qubits can be reduced to  $a + \lceil \log_2(L + 1) \rceil$ .  $\diamond$

Note that the matrix power  $A^L$  is a special case of multiplying  $L$  matrices. However, the method in Example 9.12 for encoding  $A^L$  can be highly inefficient. To see this, consider a matrix  $A$  with spectral radius

$$(9.46) \quad \rho(A) = \max \{ |\lambda| \mid \lambda \in \text{Spec}(A) \},$$

where  $\text{Spec}(A)$  denotes the set of eigenvalues of  $A$ . Suppose that  $\rho(A) < 1$ . Then there exists a constant  $C$  such that  $\sup_{L \in \mathbb{N}} \|A^L\| \leq C$ . However, it is still possible that  $\|A\| > 1$ , which means

that the block encoding subnormalization factor of  $A$  must satisfy  $\alpha \geq \|A\| > 1$ . As a result, the subnormalization factor for encoding  $A^L$  using the method in Example 9.12 would scale as  $\alpha^L$ , growing exponentially with  $L$ . This discrepancy in computing matrix powers is closely related to the challenges of solving linear differential equations. This is a topic that will be discussed in ??.

#### 9.4. Example: implementing generalized measurements

#### 9.5. Example: Quantum error correction as block encoding

#### 9.6. Query models for matrix entries

Throughout the discussion we assume  $A$  is an  $n$ -qubit, square matrix, and the max norm of  $A$  (see Definition 2.45) satisfies  $\|A\|_{\max} < 1$ .

To query the entries of a matrix, one of the most convenient form is to encode the information of the matrix as the amplitude of a known vector, e.g.,

$$(9.47) \quad O_A |0\rangle |i\rangle |j\rangle = \left( A_{ij} |0\rangle + \sqrt{1 - |A_{ij}|^2} |1\rangle \right) |i\rangle |j\rangle.$$

In other words, given  $i, j \in [N]$ ,  $O_A$  performs a controlled rotation (controlling on  $i, j$ ) on the ancilla qubit, which encodes the information in terms of amplitude of  $|0\rangle$ . We refer to Eq. (9.47) as the **amplitude oracle** or **phase oracle**.

**Example 9.13** (Construction of amplitude oracle). Assume  $\|A\|_{\max} < 1$  and  $A_{ij} \in \mathbb{R}$  for all  $i, j$ , and that we have access to a **bit oracle**

$$(9.48) \quad \tilde{O}_A |0^{d'}\rangle |i\rangle |j\rangle = |\tilde{A}_{ij}\rangle |i\rangle |j\rangle.$$

Here  $\tilde{A}_{ij}$  is a  $d'$ -bit fixed point representation of  $A_{ij}$ , and the value of  $\tilde{A}_{ij}$  is either computed on-the-fly with a quantum computer, or obtained through an external database using e.g., QRAM in Definition 5.6. Using the classical arithmetic operations, we can first convert this oracle into an oracle

$$(9.49) \quad O'_A |0^d\rangle |i\rangle |j\rangle = |\tilde{\theta}_{ij}\rangle |i\rangle |j\rangle,$$

where  $0 \leq \tilde{\theta}_{ij} < 1$ , and  $\tilde{\theta}_{ij}$  is a  $d$ -bit representation of  $\theta_{ij} = \arccos(A_{ij})/\pi$ , and with some abuse of notation we redefine  $\tilde{A}_{ij} = \cos(\pi\tilde{\theta}_{ij})$ . This step may require some additional work registers not shown here.

Now using the controlled rotation in Proposition 5.7 and Fig. 5.3, the information of  $\tilde{\theta}_{ij}$  can now be transferred to the amplitude of the ancilla qubit. We should then perform uncomputation and free the work register storing such intermediate information  $\tilde{\theta}_{ij}$ . The procedure is as follows

$$(9.50) \quad \begin{aligned} & |0\rangle \underbrace{|0^d\rangle}_{\text{work register}} |i\rangle |j\rangle \xrightarrow{I_1 \otimes O'_A} |0\rangle |\tilde{\theta}_{ij}\rangle |i\rangle |j\rangle \\ & \xrightarrow{\text{CR}} \left( \tilde{A}_{ij} |0\rangle + \sqrt{1 - |\tilde{A}_{ij}|^2} |1\rangle \right) |\tilde{\theta}_{ij}\rangle |i\rangle |j\rangle \\ & \xrightarrow{I_1 \otimes (O'_A)^{-1}} \left( \tilde{A}_{ij} |0\rangle + \sqrt{1 - |\tilde{A}_{ij}|^2} |1\rangle \right) |0^d\rangle |i\rangle |j\rangle \end{aligned}$$

After the uncomputation, the  $d$ -bit working register can be discarded, and we obtain the desired amplitude oracle of the input matrix  $A$ .  $\diamond$

**Exercise 9.4.** Construct a query oracle  $O_A$  similar to that in Eq. (9.50), when  $A_{ij} \in \mathbb{C}$  with  $\|A\|_{\max} < 1$ .

### 9.7. Block encoding of $s$ -sparse matrices

**Example 9.14** (Block encoding of a diagonal matrix). As a special case, let us consider the block encoding of a diagonal matrix, which is also a 1-sparse matrix. Since the row and column indices are the same, we may simplify the oracle Eq. (9.47) into

$$(9.51) \quad O_A |0\rangle |i\rangle = \left( A_{ii} |0\rangle + \sqrt{1 - |A_{ii}|^2} |1\rangle \right) |i\rangle.$$

Let  $U_A = O_A$ . Direct calculation shows that for any  $i, j \in [N]$ ,

$$(9.52) \quad \langle 0 | \langle i | U_A | 0 \rangle | j \rangle = A_{ii} \delta_{ij}.$$

This proves that  $U_A \in \text{BE}_{1,1}(A)$ , i.e.,  $U_A$  is a  $(1, 1)$ -block-encoding of the diagonal matrix  $A$ .  $\diamond$

**Example 9.15** (General 1-sparse matrices). In a 1-sparse matrices, there is only one nonzero entry in each row and each column of the matrix. This means that for each  $j \in [N]$ , there is a unique  $c(j) \in [N]$  such that  $A_{c(j),j} \neq 0$ , and the mapping  $c$  is a permutation. Assume that there exists a unitary  $O_c$  satisfying that

$$(9.53) \quad O_c |j\rangle = |c(j)\rangle, \quad O_c^\dagger |c(j)\rangle = |j\rangle.$$

The implementation of  $O_c$  may require the usage of some work registers that are omitted here.

We assume the matrix entry  $A_{c(j),j}$  can be queried via

$$(9.54) \quad O_A |0\rangle |j\rangle = \left( A_{c(j),j} |0\rangle + \sqrt{1 - |A_{c(j),j}|^2} |1\rangle \right) |j\rangle.$$

Now we construct  $U_A = (I \otimes O_c) O_A$ , and compute the matrix element

$$(9.55) \quad \langle 0 | \langle i | U_A | 0 \rangle | j \rangle = \langle 0 | \langle i | \left( A_{c(j),j} |0\rangle + \sqrt{1 - |A_{c(j),j}|^2} |1\rangle \right) |c(j)\rangle = A_{c(j),j} \delta_{i,c(j)}.$$

This proves that  $U_A \in \text{BE}_{1,1}(A)$ .  $\diamond$

For a general  $s$ -sparse matrix, we have  $\|A\| \leq s \|A\|_{\max}$  according to Lemma 2.46, and the explicit construction of the block encoding matrix often requires to choose the subnormalization factor  $\alpha = s \|A\|_{\max}$ . WLOG we assume each row and each column has exactly  $s$  nonzero entries (otherwise we can always treat some zero entries as nonzeros). For simplicity, let  $s = 2^s$ . For each column  $j$ , the row index for the  $\ell$ -th nonzero entry is denoted by  $c(j, \ell)$ .

**Example 9.16** (Banded matrix). In a banded matrix, we have

$$(9.56) \quad c(j, \ell) = j + \ell - \ell_0 \pmod{N},$$

for some  $\ell_0 \in \mathbb{Z}$ . The bandwidth is  $s$ . Using an adder circuit to perform the addition of  $j, \ell$  coherently, we can construct a unitary  $O_c$  such that

$$(9.57) \quad O_c |\ell\rangle |j\rangle = |\ell\rangle |c(j, \ell)\rangle.$$

Here the first register is an  $s$ -qubit register. It also holds that  $O_c^\dagger |\ell\rangle |c(j, \ell)\rangle = |\ell\rangle |j\rangle$ . This means that for each row index  $i = c(j, \ell)$ , we can recover the column index  $j$  given the value of  $\ell$ . From Eq. (9.57), we assume that the matrix entries can be queried via

$$(9.58) \quad O_A |0\rangle |\ell\rangle |j\rangle = \left( A_{c(j, \ell), j} |0\rangle + \sqrt{1 - |A_{c(j, \ell), j}|^2} |1\rangle \right) |\ell\rangle |j\rangle.$$

We then define  $D = H^{\otimes s}$  (the  $s$ -qubit Hadamard transform) satisfying

$$(9.59) \quad D |0^s\rangle = \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} |\ell\rangle.$$

We now claim that the circuit in Fig. 9.7 defines a unitary  $U_A$  that is a  $(s, s+1, 0)$ -block encoding of  $A$ .

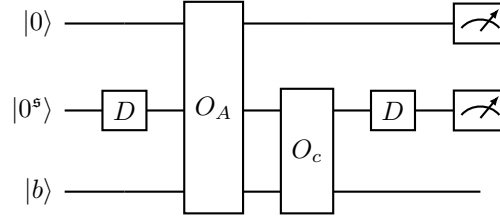


FIGURE 9.7. Quantum circuit for block encoding a banded matrix. The measurement means that to obtain a state  $\propto A|b\rangle$ , the ancilla register should all return the value 0.

◇

**Proposition 9.17.** *The circuit in Fig. 9.7 defines  $U_A \in \text{BE}_{s,s+1}(A)$ .*

PROOF. We may write

$$(9.60) \quad U_A = (I \otimes D \otimes I)(I \otimes O_c)O_A(I \otimes D \otimes I).$$

In order to compute the inner product  $\langle 0| \langle 0^s| \langle i| U_A |0\rangle |0^s\rangle |j\rangle$ , we apply  $D, O_A, O_c$  to  $|0\rangle |0^s\rangle |j\rangle$  successively as

$$(9.61) \quad \begin{aligned} |0\rangle |0^s\rangle |j\rangle &\xrightarrow{D} \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} |0\rangle |\ell\rangle |j\rangle \\ &\xrightarrow{O_A} \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} \left( A_{c(j,\ell),j} |0\rangle + \sqrt{1 - |A_{c(j,\ell),j}|^2} |1\rangle \right) |\ell\rangle |j\rangle \\ &\xrightarrow{O_c} \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} \left( A_{c(j,\ell),j} |0\rangle + \sqrt{1 - |A_{c(j,\ell),j}|^2} |1\rangle \right) |\ell\rangle |c(j,\ell)\rangle. \end{aligned}$$

Instead of multiplying the leftmost factor  $I \otimes D \otimes I$  to the last line, we apply it to  $|0\rangle |0^s\rangle |i\rangle$  first to obtain (note that  $D$  is Hermitian)

$$(9.62) \quad |0\rangle |0^s\rangle |i\rangle \xrightarrow{D} \frac{1}{\sqrt{s}} \sum_{\ell' \in [s]} |0\rangle |\ell'\rangle |i\rangle.$$

Finally, taking the inner product yields

$$(9.63) \quad \langle 0| \langle 0^s| \langle i| U_A |0\rangle |0^s\rangle |j\rangle = \frac{1}{s} \sum_{\ell} A_{c(j,\ell),j} \delta_{i,c(j,\ell)} = \frac{1}{s} A_{ij}.$$

□

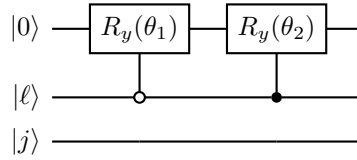
**Example 9.18.** Let us use the circuit in Fig. 9.7 to construct a block encoding of

$$(9.64) \quad A = \begin{bmatrix} \alpha_1 & \alpha_2 \\ \alpha_2 & \alpha_1 \end{bmatrix}, \quad 0 \leq \alpha_i \leq 1, \quad i = 1, 2.$$

This matrix satisfies  $\|A\|_{\max} = 1$ , and can be viewed as a 2-sparse, banded matrix. We can simply use CNOT as the  $O_c$  circuit by examining the truth table

$$\begin{array}{cc|cc} \ell & j & & \ell & c(j, \ell) \\ \hline 0 & 0 & & 0 & 0 \\ 0 & 1 & \rightarrow & 0 & 1 \\ 1 & 0 & & 1 & 1 \\ 1 & 1 & & 1 & 0 \end{array}$$

Meanwhile  $O_A$  can be implemented using controlled  $R_y(\theta_i), \theta_i = 2 \arccos(\alpha_i), i = 1, 2$ .



For example, when  $\alpha_1 = 1, \alpha_2 = 0.5$ , The resulting matrix is

$$(9.65) \quad U_A = \begin{pmatrix} 0.500 & 0.250 & 0.500 & -0.250 & 0.0 & -0.433 & 0.0 & 0.433 \\ 0.250 & 0.500 & -0.250 & 0.500 & -0.433 & 0.0 & 0.433 & 0.0 \\ 0.500 & -0.250 & 0.500 & 0.250 & 0.0 & 0.433 & 0.0 & -0.433 \\ -0.250 & 0.500 & 0.250 & 0.500 & 0.433 & 0.0 & -0.433 & 0.0 \\ 0.0 & 0.433 & 0.0 & -0.433 & 0.500 & 0.250 & 0.500 & -0.250 \\ 0.433 & 0.0 & -0.433 & 0.0 & 0.250 & 0.500 & -0.250 & 0.500 \\ 0.0 & -0.433 & 0.0 & 0.433 & 0.500 & -0.250 & 0.500 & 0.250 \\ -0.433 & 0.0 & 0.433 & 0.0 & -0.250 & 0.500 & 0.250 & 0.500 \end{pmatrix}.$$

This is a  $(2, 2)$ -block-encoding of  $A$ . ◇

**Exercise 9.5.** Construct an  $s$ -sparse matrix so that the oracle of the form Eq. (9.57) does not exist.

For more general  $s$ -sparse matrices, we need to consider a more general input model to construct its block encoding. We assume access to the following two  $(2n)$ -qubit oracles

$$(9.66) \quad \begin{aligned} O_r |\ell\rangle |i\rangle &= |r(i, \ell)\rangle |i\rangle, \\ O_c |\ell\rangle |j\rangle &= |c(j, \ell)\rangle |j\rangle. \end{aligned}$$

Here  $r(i, \ell), c(j, \ell)$  gives the  $\ell$ -th nonzero entry in the  $i$ -th row and  $j$ -th column, respectively. It should be noted that although the index  $\ell \in [s]$ , we should expand it into an  $n$ -qubit state (e.g. let  $\ell$  take the last  $s$  qubits of the  $n$ -qubit register following the binary representation of integers).

Similar to the discussion before, we need an operator  $D$  satisfying

$$(9.67) \quad D |0^n\rangle = \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} |\ell\rangle.$$

This can be implemented using Hadamard gates as

$$(9.68) \quad D = H^{\otimes s} \otimes I_{n-s}.$$

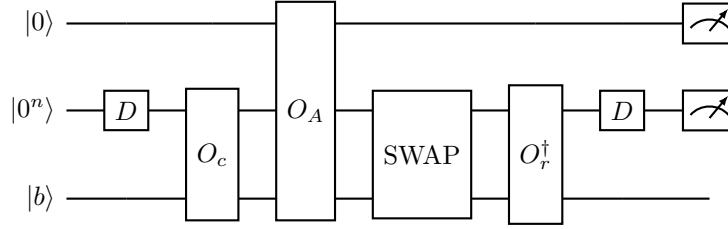


FIGURE 9.8. Quantum circuit for block encoding of general sparse matrices. The measurement means that to obtain a state  $\propto A|b\rangle$ , the ancilla register should all return the value 0.

We assume that the matrix entries are queried using the following oracle using controlled rotations

$$(9.69) \quad O_A |0\rangle |i\rangle |j\rangle = \left( A_{ij} |0\rangle + \sqrt{1 - |A_{ij}|^2} |1\rangle \right) |i\rangle |j\rangle,$$

where the rotation is controlled by both row and column indices. However, if  $A_{ij} = 0$  for some  $i, j$ , the rotation can be arbitrary, as there will be no contribution due to the usage of  $O_r, O_c$ .

**Proposition 9.19.** *Fig. 9.8 defines  $U_A \in \text{BE}_{s,n+1}(A)$ .*

PROOF. We apply the first four gate sets to the source state

$$(9.70) \quad \begin{aligned} & |0\rangle |0^n\rangle |j\rangle \\ & \xrightarrow{D, O_c, O_A} \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} \left( A_{c(j, \ell), j} |0\rangle + \sqrt{1 - |A_{c(j, \ell), j}|^2} |1\rangle \right) |c(j, \ell)\rangle |j\rangle \\ & \xrightarrow{\text{SWAP}} \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} \left( A_{c(j, \ell), j} |0\rangle + \sqrt{1 - |A_{c(j, \ell), j}|^2} |1\rangle \right) |j\rangle |c(j, \ell)\rangle. \end{aligned}$$

We then apply  $D$  and  $O_r$  to the target state

$$(9.71) \quad |0\rangle |0^n\rangle |i\rangle \xrightarrow{D, O_r} \frac{1}{\sqrt{s}} \sum_{\ell' \in [s]} |0\rangle |r(i, \ell')\rangle |i\rangle.$$

Then the inner product gives

$$(9.72) \quad \begin{aligned} \langle 0 | \langle 0^n | \langle i | U_A | 0 \rangle | 0^n \rangle | j \rangle &= \frac{1}{s} \sum_{\ell, \ell'} A_{c(j, \ell), j} \delta_{i, c(j, \ell)} \delta_{r(i, \ell'), j} \\ &= \frac{1}{s} \sum_{\ell} A_{c(j, \ell), j} \delta_{i, c(j, \ell)} = \frac{1}{s} A_{ij}. \end{aligned}$$

If  $A_{ij} \neq 0$ , then there exists a unique  $\ell$  such that  $i = c(j, \ell)$  and a unique  $\ell'$  such that  $j = r(i, \ell')$ ; if  $A_{ij} = 0$ , then the same computation gives  $\langle 0 | \langle 0^n | \langle i | U_A | 0 \rangle | 0^n \rangle | j \rangle = 0$ .  $\square$

We remark that the quantum circuit in Fig. 9.8 is essentially the construction in [GSLW18, Lemma 48], which gives a  $(s, n + 3)$ -block-encoding. The construction above slightly simplifies the procedure and saves two extra qubits (used to mark whether  $\ell \geq s$ ).

### 9.8. Hermitian block encoding

So far we have considered general  $s$ -sparse matrices. Note that if  $A$  is a Hermitian matrix, its  $(\alpha, m, \epsilon)$ -block-encoding  $U_A$  does not need to be Hermitian. Even if  $\epsilon = 0$ , we only have that the upper-left  $n$ -qubit block of  $U_A$  is Hermitian. For instance, even the block encoding of a Hermitian, diagonal matrix in Example 9.14 may not be Hermitian. On the other hand, there are cases when  $U_A = U_A^\dagger$  is a Hermitian matrix, and hence the definition:

**Definition 9.20** (Hermitian block encoding). *Let  $U_A$  be an  $(\alpha, m, \epsilon)$ -block-encoding of  $A$ . If  $U_A$  is also Hermitian, then it is called an  $(\alpha, m, \epsilon)$ -Hermitian-block-encoding of  $A$ . When  $\epsilon = 0$ , it is called an  $(\alpha, m)$ -Hermitian-block-encoding. The set of all  $(\alpha, m, \epsilon)$ -Hermitian-block-encodings of  $A$  is denoted by  $\text{HBE}_{\alpha, m}(A, \epsilon)$ , and we define  $\text{HBE}_{\alpha, m}(A) = \text{HBE}_{\alpha, m}(A, 0)$ .*

The Hermitian block encoding provides the simplest scenario of the qubitization process in ??.

Next we consider the Hermitian block encoding of an  $s$ -sparse Hermitian matrix. Since  $A$  is Hermitian, we only need one oracle to query the location of the nonzero entries

$$(9.73) \quad O_c | \ell \rangle | j \rangle = | c(j, \ell) \rangle | j \rangle .$$

Here  $c(j, \ell)$  gives the  $\ell$ -th nonzero entry in the  $j$ -th column. It can also be interpreted as the  $\ell$ -th nonzero entry in the  $j$ -th row. Again the first register needs to be interpreted as an  $n$ -qubit register. The operator  $D$  is the same as in Eq. (9.68).

Unlike all discussions before, we introduce two control qubits, and a quantum state in the computational basis takes the form  $|a\rangle |i\rangle |b\rangle |j\rangle$ , where  $a, b \in \{0, 1\}$ ,  $i, j \in [N]$ . In other words, we may view  $|a\rangle |i\rangle$  as the first register, and  $|b\rangle |j\rangle$  as the second register. The  $(n+1)$ -qubit SWAP gate is defined as

$$(9.74) \quad \text{SWAP} |a\rangle |i\rangle |b\rangle |j\rangle = |b\rangle |j\rangle |a\rangle |i\rangle .$$

To query matrix entries, we need access to the square root of  $A_{ij}$  as (note that act on the second single-qubit register)

$$(9.75) \quad O_A |i\rangle |0\rangle |j\rangle = |i\rangle \left( \sqrt{A_{ij}} |0\rangle + \sqrt{1 - |A_{ij}|} |1\rangle \right) |j\rangle .$$

Throughout we assume  $\|A\|_{\max} \leq 1$ , so that the right-hand side is normalized. The square root operation is well defined if  $A_{ij} \geq 0$  for all entries. If  $A$  has negative (or complex) entries, we first write  $A_{ij} = |A_{ij}| e^{i\theta_{ij}}$ ,  $\theta_{ij} \in [0, 2\pi)$ , and the square root is uniquely defined as  $\sqrt{A_{ij}} = \sqrt{|A_{ij}|} e^{i\theta_{ij}/2}$ .

**Proposition 9.21.** *Fig. 9.9 defines  $U_A \in \text{HBE}_{s, n+2}(A)$ .*

PROOF. Apply the first four gate sets to the source state gives

$$(9.76) \quad \begin{aligned} & |0\rangle |0^n\rangle |0\rangle |j\rangle \xrightarrow{D} \xrightarrow{O_c} \\ & \xrightarrow{O_A} \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} |0\rangle |c(j, \ell)\rangle \left( \sqrt{A_{c(j, \ell), j}} |0\rangle + \sqrt{1 - |A_{c(j, \ell), j}|} |1\rangle \right) |j\rangle \\ & \xrightarrow{\text{SWAP}} \frac{1}{\sqrt{s}} \sum_{\ell \in [s]} \left( \sqrt{A_{c(j, \ell), j}} |0\rangle + \sqrt{1 - |A_{c(j, \ell), j}|} |1\rangle \right) |j\rangle |0\rangle |c(j, \ell)\rangle \end{aligned}$$



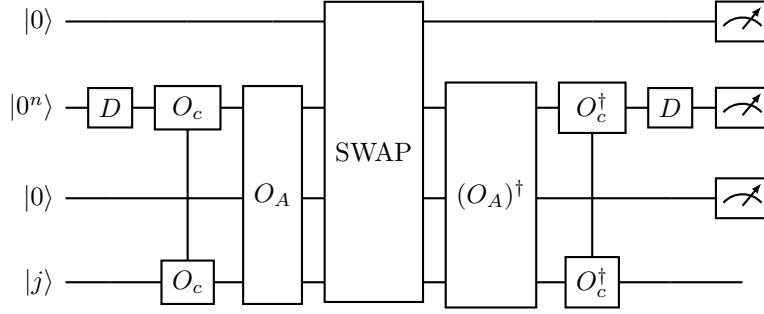


FIGURE 9.9. Quantum circuit for Hermitian block encoding of a general Hermitian matrix

Apply the last three gate sets to the target state

$$(9.77) \quad \begin{aligned} & |0\rangle |0^n\rangle |0\rangle |i\rangle \xrightarrow{D, O_c} \\ & \xrightarrow{O_A} \frac{1}{\sqrt{s}} \sum_{\ell' \in [s]} |0\rangle |c(i, \ell')\rangle \left( \sqrt{A_{c(i, \ell'), i}} |0\rangle + \sqrt{1 - |A_{c(i, \ell'), i}|} |1\rangle \right) |i\rangle \end{aligned}$$

Finally, take the inner product as

$$(9.78) \quad \begin{aligned} & \langle 0 | \langle 0^n | \langle 0 | \langle i | U_A | 0 \rangle | 0^n \rangle | 0 \rangle | j \rangle \\ &= \frac{1}{s} \sum_{\ell, \ell'} \sqrt{A_{c(j, \ell), j}} \sqrt{A_{c(i, \ell'), i}^*} \delta_{i, c(j, \ell)} \delta_{c(i, \ell'), j} \\ &= \frac{1}{s} \sqrt{A_{ij}} \sqrt{A_{ji}^*} = \frac{1}{s} (\sqrt{A_{ij}})^2 = \frac{1}{s} A_{ij}. \end{aligned}$$

In this equality, we have used that  $A$  is Hermitian:  $A_{ij} = A_{ji}^*$ , and there exists a unique  $\ell$  such that  $i = c(j, \ell)$ , as well as a unique  $\ell'$  such that  $j = c(i, \ell')$  when  $A_{ij}$  is nonzero.  $\square$

**Exercise 9.6.** Let  $A \in \mathbb{C}^{N \times N}$  ( $N = 2^n$ ) be a Hermitian matrix with entries on the complex unit circle  $A_{ij} = e^{i\theta_{ij}}$ ,  $\theta_{ij} \in [0, 2\pi)$ , which can be accessed via a  $2n$  qubit unitary  $V \in \mathbb{C}^{N^2 \times N^2}$  such that

$$V |0^n\rangle |j\rangle = \frac{1}{\sqrt{N}} \sum_{i \in [N]} e^{i\theta_{ij}/2} |i\rangle |j\rangle, \quad j \in [N].$$

Use  $V$  to implement a block encoding  $U$  of  $A$  with  $n$  ancilla qubits. What is the subnormalization factor  $\alpha$  for this block encoding?

### Notes and further reading

The mathematical idea underlying block encodings is a form of unitary dilation: linear maps that are not themselves unitary can often be realized as a sub-block of a larger unitary acting on an extended space. In quantum information, this viewpoint is closely related to dilation theorems for completely positive maps. In quantum algorithms, the block-encoding terminology (together with explicit bookkeeping of the subnormalization factor and approximation error) was systematized as part of the modern polynomial-transformation framework; see [GSLW19].

The linear combination of unitaries (LCU) primitive used here originates in the Hamiltonian simulation algorithm [CW12, BCC<sup>+</sup>14]. In particular, the sparse-matrix block-encoding constructions in this chapter are closely aligned with the query models developed for sparse Hamiltonian simulation (see, e.g., [BACS07]) and with the block-encoding-based linear-systems framework (see, e.g., [CKS17], which can be directly connected to the quantum circuit for Hermitian block encoding in Fig. 9.9). The connection between block encodings and quantum walks is mediated by the fact that many walk operators are themselves natural block encodings; see Szegedy’s quantization of Markov chains [Sze04] for an early and influential formulation, which will be discussed in detail in ???. The RACBEM input model for pseudorandom nonunitary matrices was introduced in [DL21].

## Bibliography

- [AA11] Scott Aaronson and Alex Arkhipov. The computational complexity of linear optics. In **Proceedings of the forty-third annual ACM symposium on Theory of computing**, pages 333–342, 2011.
- [Aar14] Scott Aaronson. Quantum machine learning algorithms: Read the fine print. **Nat. Phys.**, page 5, 2014.
- [ABO97] Dorit Aharonov and Michael Ben-Or. Fault-tolerant quantum computation with constant error. In **Proceedings of the twenty-ninth annual ACM symposium on Theory of computing**, pages 176–188, 1997.
- [BACS07] Dominic W Berry, Graeme Ahokas, Richard Cleve, and Barry C Sanders. Efficient quantum algorithms for simulating sparse hamiltonians. **Commun. Math. Phys.**, 270:359–371, 2007.
- [BBC<sup>+</sup>95] Adriano Barenco, Charles H Bennett, Richard Cleve, David P DiVincenzo, Norman Margolus, Peter Shor, Tycho Sleator, John A Smolin, and Harald Weinfurter. Elementary gates for quantum computation. **Phys. Rev. A**, 52:3457, 1995.
- [BBK<sup>+</sup>23] Ryan Babbush, Dominic W Berry, Robin Kothari, Rolando D Somma, and Nathan Wiebe. Exponential quantum speedup in simulating coupled classical oscillators. **Phys. Rev. X**, 13:041041, 2023.
- [BCC<sup>+</sup>14] Dominic W Berry, Andrew M Childs, Richard Cleve, Robin Kothari, and Rolando D Somma. Exponential improvement in precision for simulating sparse Hamiltonians. In **Proceedings of the forty-sixth annual ACM symposium on Theory of computing**, pages 283–292, 2014.
- [Ben87] Charles H Bennett. Demons, engines and the second law. **Scientific American**, 257:108–117, 1987.
- [Bha97] Rajendra Bhatia. **Matrix Analysis**, volume 169. Springer, 1997.
- [Bri98] Matthew Edward Briggs. **An introduction to the general number field sieve**. PhD thesis, Virginia Tech, 1998.
- [CCHL22] Sitan Chen, Jordan Cotler, Hsin-Yuan Huang, and Jerry Li. Exponential separations between learning with and without quantum memory. In **2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)**, pages 574–585. IEEE, 2022.
- [Cho75] Man-Duen Choi. Completely positive linear maps on complex matrices. **Linear algebra and its applications**, 10(3):285–290, 1975.
- [CKS17] Andrew M. Childs, Robin Kothari, and Rolando D. Somma. Quantum algorithm for systems of linear equations with exponentially improved dependence on precision. **SIAM J. Comput.**, 46:1920–1950, 2017.
- [CW12] Andrew M. Childs and Nathan Wiebe. Hamiltonian simulation using linear combinations of unitary operations. **Quantum Information and Computation**, 12:901–924,

- 2012.
- [DL21] Yulong Dong and Lin Lin. Random circuit block-encoded matrix and a proposal of quantum linpack benchmark. **Phys. Rev. A**, 103:062412, 2021.
- [Fey82] Richard P Feynman. Simulating physics with computers. **Int. J. Theor. Phys**, 21, 1982.
- [FVDG02] Christopher A Fuchs and Jeroen Van De Graaf. Cryptographic distinguishability measures for quantum-mechanical states. **IEEE T. Inform. Theory**, 45:1216–1227, 2002.
- [Gid18] Craig Gidney. Halving the cost of quantum addition. **Quantum**, 2:74, 2018.
- [GLM08] Vittorio Giovannetti, Seth Lloyd, and Lorenzo Maccone. Quantum random access memory. **Physical review letters**, 100(16):160501, 2008.
- [GSLW18] András Gilyén, Yuan Su, Guang Hao Low, and Nathan Wiebe. Quantum singular value transformation and beyond: exponential improvements for quantum matrix arithmetics. **arXiv:1806.01838**, 2018.
- [GSLW19] András Gilyén, Yuan Su, Guang Hao Low, and Nathan Wiebe. Quantum singular value transformation and beyond: exponential improvements for quantum matrix arithmetics. In **Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing**, pages 193–204, 2019.
- [GVL13] G. H. Golub and C. F. Van Loan. **Matrix computations**. Johns Hopkins Univ. Press, 2013.
- [HBC<sup>+</sup>22] Hsin-Yuan Huang, Michael Broughton, Jordan Cotler, Sitan Chen, Jerry Li, Masoud Mohseni, Hartmut Neven, Ryan Babbush, Richard Kueng, John Preskill, et al. Quantum advantage in learning from experiments. **Science**, 376(6598):1182–1186, 2022.
- [Hel69] Carl W Helstrom. Quantum detection and estimation theory. **Journal of Statistical Physics**, 1(2):231–252, 1969.
- [Hig08] N. Higham. **Functions of matrices: theory and computation**, volume 104. SIAM, 2008.
- [Jam72] Andrzej Jamiolkowski. Linear transformations which preserve trace and positive semidefiniteness of operators. **Reports on mathematical physics**, 3(4):275–278, 1972.
- [JSW<sup>+</sup>25] Stephen P Jordan, Noah Shutty, Mary Wootters, Adam Zalcman, Alexander Schimhuber, Robbie King, Sergei V Isakov, Tanuj Khattar, and Ryan Babbush. Optimization by decoded quantum interferometry. **Nature**, 646(8086):831–836, 2025.
- [KBDW83] Karl Kraus, Arno Böhm, John D Dollard, and WH Wootters. **States, effects, and operations fundamental notions of quantum theory: Lectures in mathematical physics at the university of Texas at Austin**. Springer, 1983.
- [Kit97] A Yu Kitaev. Quantum computations: algorithms and error correction. **Russian Mathematical Surveys**, 52(6):1191, 1997.
- [Lan61] Rolf Landauer. Irreversibility and heat generation in the computing process. **IBM journal of research and development**, 5:183–191, 1961.
- [LC17] Guang Hao Low and Isaac L. Chuang. Optimal hamiltonian simulation by quantum signal processing. **Phys. Rev. Lett.**, 118:010501, 2017.
- [Llo96] Seth Lloyd. Universal quantum simulators. **Science**, pages 1073–1078, 1996.
- [Man80] Yu I Manin. Vychislimoe i nevychislimoe (computable and noncomputable), moscow: Sov, 1980.
- [NC00] Michael A Nielsen and Isaac Chuang. Quantum computation and quantum information, 2000.

- [RP11] Eleanor G Rieffel and Wolfgang H Polak. **Quantum computing: A gentle introduction**. MIT Press, 2011.
- [Sho94] Peter W Shor. Algorithms for quantum computation: discrete logarithms and factoring. In **Proceedings 35th annual symposium on foundations of computer science**, pages 124–134. Ieee, 1994.
- [Sho99] Peter W Shor. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. **SIAM review**, 41:303–332, 1999.
- [Sti55] W Forrest Stinespring. Positive functions on  $c^*$ -algebras. **Proceedings of the American Mathematical Society**, 6(2):211–216, 1955.
- [Sze04] Mario Szegedy. Quantum speed-up of markov chain based algorithms. In **45th Annual IEEE symposium on foundations of computer science**, pages 32–41, 2004.
- [TM71] Myron Tribus and Edward C McIrvine. Energy and information. **Scientific American**, 225:179–190, 1971.
- [Tur36] Alan Turing. On computable numbers, with an application to the entscheidungsproblem. **J. Math**, 58:5, 1936.
- [Uhl76] Armin Uhlmann. The transition probability in the state space of a  $c^*$  algebra. **Reports on Mathematical Physics**, 9(2):273–279, 1976.
- [VN93] John Von Neumann. First draft of a report on the edvac. **IEEE Ann. Hist. Comput.**, 15:27–75, 1993.
- [Wat18] John Watrous. **The theory of quantum information**. Cambridge Univ. Pr., 2018.



# Index

- T gate, 25
- $s$ -sparse matrix, 46
  
- adjoint map, 83
- amplitude oracle, 123
- angle, 80
- asymptotic notations, 22
  
- Baker–Campbell–Hausdorff formula, 26
- Bell state, 36
- bit oracle, 123
- Bloch sphere, 24
- Block encoding, 115
- bosonic operator, 48
- bra vector, 23
  
- canonical anticommutation relations, 48
- canonical commutation relations, 48
- Choi matrix, 66
- Choi–Jamiolkowski isomorphism, 66
- classical channel, 65
- classical ensemble, 34
- classical state, 40, 65
- Clifford group, 31
- CNOT gate, 30
- completely positive, 63, 67
- complex projective space, 70
- computational basis, 40
- controlled unitary, 30
  
- density matrix, 34
- density operator, 34
- dephasing channel, 65
- distance, 69
  
- entanglement, 8
- equivalence relation, 70
- Extended Church–Turing Thesis, 8
  
- fermionic operator, 48
- fidelity, 79
- fixed point representation, 99
  
- generalized measurement, 36
- global phase invariant distance, 70
  
- Hadamard gate, 25
- Hadamard test circuit, 40
- Hamiltonian, 24
- Hermitian matrix, 21
- Hilbert space, 23
  
- identity channel, 62
- induced total variation distance, 73
- induced trace distance, 89
- induced trace norm, 82
- induced vector 2-norm, 22
  
- Jordan–Wigner transformation, 47
  
- ket vector, 23
- ketbra notation, 23
- Kraus form, 65
  
- linear combination of unitaries, 116
  
- Majorana operator, 47
- matrix exponential, 26
- matrix function, 26
- matrix norm
  - Schatten 1-norm, 22, 74
  - Schatten 2-norm, 74
  - Schatten  $\infty$ -norm, 22, 74
  - Schatten  $p$ -norm, 22, 74
  - trace norm, 22
- max norm, 46
- metric, 69
- mixed state, 34
- mixed state, maximally, 34
  
- Naimark’s dilation theorem, 37
- no-cloning theorem, 42
- normal matrix, 21
  - spectral theorem of, 25
- operator exponential, 26

- operator norm, 22
- operator sum representation, 65
- oracle, 97
- partial application of operators, 33
- partial inner product, 32
- partial trace, 35
- Pauli group, 31
- Pauli matrices, 25
- phase gate, 25
- phase oracle, 123
- positive operator, 22, 63
- positive operator-valued measure, 36
- prepare oracle, 117
- principle of deferred measurement, 45
- principle of implicit measurements, 45
- probabilistic state, 65
- probability distribution, 59
- product states, 28
- projective measurement, 35
- projective unitary group, 71
- pure state, 34
- quantum advantage, 10
- quantum channel, 62, 64
- quantum circuit, 38
- quantum gate, 25
- quantum learning theory, 11
- quantum measurements, 27
- quantum observable, 27
- quantum random access memory, 99
- quantum register, 40
- quantum speedup, 10
- qubits, 22
- random circuit based block-encoded matrix, 116
- real dimension, 70
- reduced density operators, 36
- Reversible computation, 95
- select oracle, 117
- singular value decomposition, 26
- state vector, 23
- superoperator, 62
- SWAP gate, 39
- SWAP test, 41
- tensor product
  - linear operator, 29
  - superoperator, 62
  - vector, 28
- Toffoli gate, 39
- total variation distance, 72
- trace distance, 77
- trace norm, 74
- trace preserving, 63
- transition matrix, 61
- uncomputation, 96
- unitary channel, 64
- vector 2-norm, 21