
Extracting Classical Information from Quantum States: Fundamental Limits, Adaptive and Finite-Length Measurements

by

Hye Won Chung

B.S., Electrical Engineering and Computer Science,
Korea Advanced Institute of Science and Technology, 2007

S.M., Electrical Engineering and Computer Science,
Massachusetts Institute of Technology, 2009

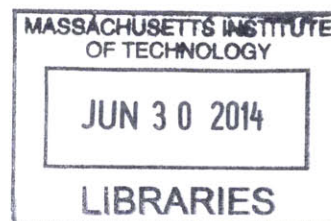
Submitted to the Department of Electrical Engineering and Computer Science
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy
in Electrical Engineering and Computer Science
at the Massachusetts Institute of Technology

June 2014

© 2014 Massachusetts Institute of Technology
All Rights Reserved.

ARCHIVES



Signature redacted

Author: _____

Department of Electrical Engineering and Computer Science
May 19, 2014

Signature redacted

Certified by: _____

Professor Lizhong Zheng
Thesis Supervisor

Signature redacted

Accepted by: _____

Professor Leslie A. Kolodziejcki
Chair, Department Committee on Graduate Students

Extracting Classical Information from Quantum States: Fundamental Limits, Adaptive and Finite-Length Measurements

by Hye Won Chung

Submitted to the Department of Electrical Engineering and Computer Science
on May 19, 2014, in partial fulfilment of the requirements for the degree of
Doctor of Philosophy in Electrical Engineering and Computer Science

Abstract

Motivated by the increasing demand for more powerful and efficient optical communication systems, quantum mechanics of information processing has become the key element in determining the fundamental limits of physical channels, and in designing quantum communication systems that approach those fundamental limits. To achieve higher data rates over quantum optical channels, we need to efficiently extract classical information from quantum states. However, peculiar properties of quantum states, such as the no-cloning theorem and the non-reversible measurement process, provide new challenges in the measurement of quantum states; in quantum information science, there is no concept analogous to sufficient statistics in classical information science. Therefore, to extract as much information as possible from quantum states, it is important to choose the right measurement process. In this thesis, we investigate the fundamental question of how to design the measurement process to efficiently extract information from quantum states.

First, we consider *adaptive measurement*, with which we measure each received quantum state one at a time, and then update the next measurement process based on the previous observations. We show that for binary hypothesis testing between two ideal laser light pulses, if we update the adaptive measurement to maximize the communication efficiency at each instant, based on recursively updated knowledge of the receiver, then we can achieve the theoretical lower bound of the detection error probability. Using this viewpoint, we give a natural generalization of the adaptive measurement to general M -ary hypothesis testing problems. We also analyze the information capacity with adaptive measurement, and compare the result with that for direct detection receivers and the ultimate capacity of quantum channels (the Holevo limit).

We also investigate finite-blocklength joint receivers. The ultimate capacity of quantum channels is calculated under the assumption that an infinite number of quantum states can be collectively measured in one shot. However, this assumption becomes the primary barrier that prevents practical implementations of capacity-achieving joint de-

tection receivers. The maximum number of classical information bits extracted per use of the quantum channel strictly increases with the number of channel outputs jointly measured at the receiver. This phenomenon is called strict superadditivity, and it has been thought of as a unique property that can be observed only in quantum channels, but not in classical discrete memoryless channels (DMCs). In this thesis, we introduce a new aspect of understanding strict superadditivity by comparing the performance of concatenated coding over quantum channels and classical DMCs, for a fixed inner code length. We show that the strict superadditivity in information rate occurs due to a loss of information from hard-decisions at the finite blocklength. We also find a lower bound on the maximum achievable information rate as a function of the length of the quantum measurement.

The analysis and new insights into the measurement process of quantum states that we develop in this thesis can be used to improve not only current quantum optical communication systems, but also classical information processing, where the data is too big to be handled with sufficient statistics. Our work would help develop new concepts of efficient statistics that provide systematic ways to choose useful information among big data while discarding the rest.

Thesis Supervisor: Professor Lizhong Zheng

Title: Professor of Electrical Engineering and Computer Science

Acknowledgments

I would like to thank my advisor Prof. Lizhong Zheng, who led me into this exciting area of information theory and quantum communications, and has advised my research throughout the past five years. His enthusiasm to understand quantum communications from information theorists' perspective motivated a great part of my research in this thesis. His insightful comments not only led me to my destination but also helped me learn many things along the way. He has been a great teacher and advisor, who is always willing to be available to discuss general perspectives of doing research as well as interesting ideas and viewpoints to understand specific problems in depth. In the future, if I have the privilege to advise students, I will try to be such a devoted advisor.

I would like to express my gratitude to my thesis committee. In the fourth year of my graduate studies, when I decided to venture into the field of quantum optical communications, Prof. Jeffrey Shapiro kindly taught me quantum optical communications from the basic principles to the most exciting research directions in this field. I deeply appreciate his guidance. The intuitions and ideas developed for quantum detection problems in this thesis are greatly inspired by Prof. Gregory Wornell's teaching. I heartily thank Prof. Gregory Wornell for his thoughtful suggestions and advice.

I am grateful to Prof. Anantha Chandrakasan, who advised my master's thesis and gave warm support when I first came to MIT. I am fortunate to have had the chance to TA with Prof. Polina Golland for Inference and Information class. I thank her for teaching me the right ways to deliver knowledge and inspire students during the class. I would like to thank Prof. Sae-Young Chung at KAIST for his continuous advice and encouragement from my undergraduate years. I am very thankful to my collaborator Dr. Saikat Guha for providing me with many insights and references.

The Claude E. Shannon Communication and Network Group members provided me motivation and inspiration during my PhD. I really enjoyed sharing offices with Mina Karzand, Fabián Kozynsky, Anuran Makur, Chung Chan, Shao-Lun Huang and Baris Nakiboglu, discussing many interesting ideas, both academic and others. I have benefited from interactions with great colleagues from RLE and EECS. Special mention goes to Vincent Tan, Da Wang, John Sun, Gauri Joshi, Ligong Wang and Qing He.

I am thankful for financial support by the Kwanjeong Educational Foundation. The work in this thesis has also been supported by the DARPA Information in a Photon (InPho) program under contract number HR0011-10-C-0159.

I am extremely grateful to my parents, Gigon Chung and Seungim Baek, my sister, Hye Jin Chung, and my husband, Ji Oon Lee, for their unconditional love and support. This thesis is dedicated to my family. Lastly, I thank God for His unfailing love and everything He has done in my life.

Contents

List of Figures	9
1 Introduction	11
1.1 Quantum Information Theory: Promises and Challenges	11
1.2 Quantum Optical Communication System	14
1.3 Adaptive Measurement	20
1.4 Joint Measurement: Superadditivity of Quantum Channel Coding Rate	23
1.5 Thesis Outline	25
2 Background	27
2.1 Quantum Hypothesis Testing	27
2.2 Capacity of Classical-Quantum Channel: Holevo capacity	33
3 Adaptive Measurements	39
3.1 Introduction	39
3.2 Preliminaries	41
3.3 Multiple-Copy States Discrimination with Adaptive Measurements . . .	43
3.4 Optimal Adaptive Measurements for Binary Hypothesis Testing	48
3.4.1 Necessary and Sufficient Conditions for Adaptive Measurement .	49
3.4.2 Bayesian Updating Rules	51
3.4.3 Dolinar Receiver	53
3.5 Conclusion	56
3.A Proof of Lemma 3.4	57
3.B Proof of Lemma 3.5	64
4 Capacity of Coherent Detection	69
4.1 Introduction: Detection of Optical Signals	69
4.2 Binary Hypothesis Testing	71
4.3 Generalization to M-ary Hypothesis Testing	75
4.4 Coded Transmissions and Capacity Results	78
4.5 Conclusion	83

4.A	Proof of Lemma 4.1	84
4.B	Proof of Lemma 4.2	85
4.C	Proof of Lemma 4.3	88
4.D	Proof of Theorem 4.4	89
5	Quantum Channel Coding Rate with Finite Blocklength Measurements	99
5.1	Background	99
5.2	Introduction and Problem Statement	101
5.3	Strict Superadditivity of C_N	106
5.4	Lower Bound on C_N	111
5.5	Proof of Theorem 5.2	115
5.5.1	Upper and Lower Bounds on the Average Probability of Error	116
5.5.2	Equierror Superchannel	118
5.6	Interpretation of Superadditivity: Classical DMC vs. Quantum Channel	120
5.6.1	A Unifying Framework to Explain Superadditivity of C_N	121
5.6.2	An Approximation of the Lower Bound on C_N	125
5.7	Conclusion	128
5.A	Proof of Lemma 5.1	129
5.B	Proof of Corollary 5.3	133
5.C	Proof of Lemma 5.4	137
6	Conclusion	143
6.1	Summary of Main Contributions	143
6.2	Suggestions for Future Research	145
6.2.1	Adaptive Measurements for M -ary Hypothesis Testing	145
6.2.2	Quantifying the Efficiency of Measurement Process: Time-Varying Metrics	147
6.2.3	Adaptive Sampling/Querying for Classical Inference Problems	148
6.2.4	Finite Blocklength Joint Measurement: Converse Bounds	149
	Bibliography	151

List of Figures

1.1	Block diagram of optical communication system	14
1.2	Concatenated coding over a classical-quantum channel	24
3.1	Coherent receiver using local feedback signal	54
4.1	Coherent Receiver Using Local Feedback Signal	70
4.2	Effective binary channel between the input hypotheses and the observation over a Δ period of time	72
4.3	An example of the control signal that achieves the minimum probability of error.	75
4.4	Empirical average of detection error probability (after 10,000 runs) for ternary hypothesis testing, using control signals that minimize the average Rényi α -entropy for different values of α ; Ternary inputs $\{ 5\rangle, -6\rangle, 3\rangle\}$ are used with prior probabilities $p = \{0.5, 0.4, 0.1\}$	78
5.1	Concatenated coding over a classical-quantum channel	103
5.2	A lower bound of photon information efficiency of the BPSK channel, $C_N/(N\mathcal{E})$, at $\mathcal{E} = 0.01$ for the finite blocklength N	114
5.3	Concatenated coding over a classical DMC	121

Introduction

■ 1.1 Quantum Information Theory: Promises and Challenges

QUANTUM information science studies the theory of communication and computation at the most fundamental physical level. Quantum communication systems transmit information in individual photons. As the demand for more powerful and efficient communication systems evolves, the quantum mechanics of information processing becomes the key element in determining the fundamental limits of physical channels such as the optical fiber or free-space optical channels, and in designing the quantum communication systems that approach the fundamental limits. For example, to help maximize the efficiency of optical communication, new techniques have been adapted that utilize a quantum mechanical phenomenon known as *entanglement*, which is a bizarre form of correlation that can be explained only by quantum mechanics [3].

One of the central questions in quantum information theory is: How many classical bits per channel use can be reliably communicated over quantum channels? In 1973, Holevo first derived an upper bound (Holevo bound) on the capacity of quantum channels to transmit classical information [22]. In [20, 23], it was shown that the Holevo bound is in principle also an achievable information rate if we can access the quantum channel infinitely many times and can implement a receiver that makes joint (collective) measurements over the infinite-length codeword blocks. The capacity of *bosonic channels*, a single electromagnetic field mode, subject to the constraint of average photon number per channel use, was first derived in the absence of noise and loss in [50]. When the average photon number per channel use per mode is \mathcal{E} , the maximum number of

information bits that can be reliably communicated through the bosonic channel is

$$C(\mathcal{E}) = (\mathcal{E} + 1) \log(\mathcal{E} + 1) - \mathcal{E} \log \mathcal{E} \text{ [nats/symbol]}. \quad (1.1)$$

In [14], this result was generalized for the purely lossy case. The capacity of bosonic channels is achievable by transmitting N -encoded input *coherent states*, i.e., ideal (single-mode) laser light of duration T -seconds, and then by jointly measuring the N -encoded coherent states, i.e., by collectively measuring the laser light of an entire duration of $N \cdot T$ -seconds, as $N \rightarrow \infty$. This is again an asymptotic result, assuming that we can store the long sequence of received coherent states until we receive the very last quantum state, and then we can jointly process those N -quantum states.

Even though the capacity of quantum channels is proven in theory, there has been no practical design of a quantum communication system that achieves the capacity of quantum channels. With current quantum optical devices, we fall short of achieving that capacity. In particular, the assumption that an arbitrarily large number of quantum states can be jointly measured is the primary barrier that prohibits practical implementations of capacity-achieving joint detection receivers (JDRs)—especially in the context of optical communication.

Current quantum optical technologies cannot store the large number of quantum states without perturbation. Moreover, in general, the complexity to implement joint detection receivers that collectively measure length- N quantum states increases exponentially with N . When considering the complexity of quantum receivers, we can ask questions as follows: How does the maximum achievable information rate increase as the length of quantum measurements increases? Or, more fundamentally, why do we need to detect a large number of quantum states *together* to achieve a higher information rate from quantum channels?

In Chapter 5 of this thesis, we will answer the first question. Even before that, the second question can be answered intuitively from *unique* properties of quantum mechanics. First, when we receive an unknown quantum state, we cannot copy it; this

property is called *no-cloning theorem* [31]. Second, once we measure a quantum state and observe its output, the quantum state is perturbed, and we cannot restore the original quantum state, or reverse the measurement process.

These two properties of quantum mechanics highly restrict our measurement process and make a great difference from the classical information processing: we have one and only one opportunity to measure the originally received quantum states, and we cannot preserve all the information in the received quantum states after we measure it. That is to say, in quantum information science, there is no concept analogous to *sufficient statistics* in classical information science. When we measure a quantum state, we can only observe *partial* information about the quantum state, and we lose the rest of the information that was originally encoded in the quantum state.

Therefore, to extract as much information as possible from quantum states, it is important to choose the right measurement. When we measure a large number of quantum states together, a much more general form of quantum measurements can be allowed, compared to the case when we measure each received quantum state one at a time. Allowing more general sets of quantum measurements may increase the amount of information that can be extracted per received state, so that we can achieve a higher information rate with a joint detection receiver that can measure a large number of quantum states at one time.

In this thesis, we investigate the fundamental question of how to design the measurement process to efficiently extract information from quantum states when the possible types of measurements are restricted to particular sets of *practically implementable* quantum receivers. For example, we consider a *coherent receiver*, which adds a feedback control signal to the received quantum state and then detects the merged signal by a photon counter. The feedback control signal that is added to the currently received quantum state can be designed based on the photon arrival histories of the previously arrived and detected quantum states. For such a fixed structure of quantum receiver, what is the optimal way to choose the feedback control signal to maximize the efficiency

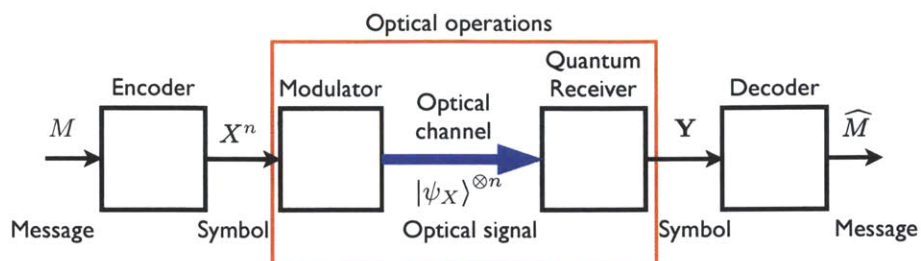


Figure 1.1. Block diagram of optical communication system

of the measurement process, by using the information extracted from the previous states to help extract new information from the next quantum state?

We will consider quantum detection/communication problems with such adaptive measurements as well as finite block-length joint measurements, and analyze the maximum achievable performance with those restricted sets of measurements. By highlighting and contrasting different properties of measurement types and their performance bounds, we will reveal what specific properties of quantum measurements indeed make a difference in the efficiency of extracting information content from quantum channels. When quantum channels paired with a certain set of measurements function analogously to a specific classical communication system, we will also use the analogy to explain the optimal ways to extract information from quantum states with that specific type of measurements.

■ 1.2 Quantum Optical Communication System

In this section, we will present how classical information bits can be communicated over physical quantum channels. Fig. 1.1 illustrates an optical communication system. The optical communication system is composed of five separate functional parts, which include an encoder, modulator, optical channel, quantum receiver, and decoder. The message we want to convey through n -uses of a quantum channel is denoted as $M \in \{1, \dots, e^{nR}\}$, where the rate of communication is R . The encoder $f(\cdot)$ maps the message M into a length n -coded input symbols (codeword) $X^n \in \mathcal{X}^n$, i.e., $f : \{1, \dots, e^{nR}\} \rightarrow$

\mathcal{X}^n .

The modulator then maps each input symbol $X \in \mathcal{X}$ into a quantum state $|\phi_X\rangle$ of optical fields, which can be transmitted through an optical medium (optical channel). The optical channel conveys the optical states $\{|\phi_X\rangle\}$, $X \in \mathcal{X}$ and maps them to a possibly different set of states, $\{|\psi_X\rangle\} \in \mathcal{H}$, $X \in \mathcal{X}$, in a Hilbert space \mathcal{H} . For the message m , whose codeword is $f(m) = x_1^n(m) = (x_1(m), x_2(m), \dots, x_n(m))$, the transmitter sends $|\phi_{x_1(m)}\rangle, |\phi_{x_2(m)}\rangle, \dots, |\phi_{x_n(m)}\rangle$ through the optical channel, and the outputs of the quantum channel become n -quantum states, $|\psi_{x_1(m)}\rangle, |\psi_{x_2(m)}\rangle, \dots, |\psi_{x_n(m)}\rangle$. The received n -quantum states for the message m can also be written as a (tensor) product state, $|\Psi_m\rangle := |\psi_{x_1(m)}\rangle \otimes |\psi_{x_2(m)}\rangle \otimes \dots \otimes |\psi_{x_n(m)}\rangle \in \mathcal{H}^{\otimes n}$.

The quantum receiver can measure the received quantum states either one at a time and generate outputs $Y_i \in \mathcal{Y}$, $i = 1, \dots, n$, or can collectively measure the length- n product state (at one shot) and generate a super-symbol $\mathbf{Y} \in \mathcal{Y}$. We can think that $\mathbf{Y} = (Y_1, \dots, Y_N)$ and $\mathcal{Y} = \mathcal{Y}^N$ for the case when the receiver detects each state one at a time. After we generate the output symbol $\mathbf{Y} \in \mathcal{Y}$ from the quantum receiver, we can decode the original message (with some probability of error) by a decoding map $g: \mathcal{Y} \rightarrow \{1, \dots, e^{nR}\}$. It is important to note that in the optical communication system illustrated in Fig. 1.1, the only parts that deal with optical signals are inside the red box, which includes the modulator, optical channel, and quantum receiver, whereas the rest of the parts are just electrical. Therefore, we can treat the input and output symbols X and \mathbf{Y} in classical ways, i.e., we can compare \mathbf{Y} with each possible X^n as many times as we want, while we have only one opportunity to measure the originally received quantum states.

We will explain the operation of the quantum receiver in more detail. The operation of the quantum receiver can be mathematically described by a Positive Operator-Valued Measure (POVM), $\{\Pi_y\}$, $y \in \mathcal{Y}$, that satisfies

$$\Pi_y \geq 0, y \in \mathcal{Y} \quad \text{and} \quad \sum_{y \in \mathcal{Y}} \Pi_y = \mathbb{1} \text{ in } \mathcal{H}^{\otimes n}, \quad (1.2)$$

where $\mathbb{1}$ is the identity operator. When the length- n product state for the message m , $|\Psi_m\rangle$, is measured by a POVM $\{\Pi_y\}$, $y \in \mathcal{Y}$, we observe y with probability

$$p(y|x_1^n(m)) = \text{Tr}(\Pi_y|\Psi_m\rangle\langle\Psi_m|) = \langle\Psi_m|\Pi_y|\Psi_m\rangle. \quad (1.3)$$

Note that when we measure each received symbol one at a time with $\{\Pi_y\}$, $y \in \mathcal{Y}$, which operates on \mathcal{H} , the probability of observing y_1^n is decomposed into n distributions, i.e.,

$$p(y_1^n|x_1^n(m)) = \prod_{i=1}^n p(y_i|x_i(m)) = \prod_{i=1}^n \langle\psi_{x_i(m)}|\Pi_{y_i}|\psi_{x_i(m)}\rangle. \quad (1.4)$$

When we fix the measurement $\{\Pi_y\}$, $y \in \mathcal{Y}$, for every use of the channel, or more generally, design POVMs for the i -th symbol $\{\Pi_y^{(i)}\}$ so as not to depend on the previous observations, the resulting channel is memoryless.

Note that we can generate much more general channel distributions, i.e., conditional probability distributions of output symbols given input symbols, when the whole sequence is *collectively* measured, as in (1.3), than in the case when each symbol is measured one at a time so that the resulting distribution is decomposed into n distributions as in (1.4).

The communication efficiency (maximum achievable information rate) as well as complexity to implement quantum receivers highly depends on the number of quantum states that can be collectively measured at the quantum receiver. There exists a trade-off between the complexity of quantum receivers, which increases exponentially with the length of joint measurements, and the maximum achievable information rate with fixed blocklengths of joint measurements. To illustrate this phenomenon more rigorously, and to understand how the optical communication happens in real physics, we present an example of an optical communication system in the rest of this section.

As one of the most common examples of optical communication systems in practice, consider a free-space (i.e., vacuum) optical communication link with coherent state inputs, i.e., ideal (single electromagnetic field mode) laser light. The modulator maps

each input symbol $X \in \mathcal{X}$ into a narrow-band optical pulse of duration T -seconds. For example, when the input symbol is binary, i.e., $X \in \{0, 1\}$, the modulator maps each symbol $X \in \{0, 1\}$ into a base-band waveform $\{S_0(t), S_1(t)\}$, $t \in [0, T)$, respectively. For the length- n codeword of message m , i.e., $f(m) = (x_1(m), x_2(m), \dots, x_n(m))$, the signal waveform can be written as

$$S(t) = \sum_{i=1}^n S_{x_i(m)}(t - i \cdot T), \quad t \in [0, nT). \quad (1.5)$$

This signal waveform is shifted up to a fixed optical angular frequency ω_s , and the modulated signal is sent through the optical channel. The transmitted optical wave (input to the optical channel) becomes

$$E_t(t) = S(t) \exp(i\omega_s t), \quad t \in [0, nT). \quad (1.6)$$

When this optical wave travels through a free-space optical communication link, the wave arriving at the receiver, denoted as $E_r(t)$, becomes an attenuated version of $E_t(t)$, and can be written as

$$E_r(t) = \sqrt{\eta} S(t) \exp(i\omega_s t), \quad t \in [0, nT), \quad (1.7)$$

where $\eta \in [0, 1]$ is an attenuation parameter depending on the distance of the communication link. Here, we have neglected the propagation delay.

For such an electromagnetic wave, the average photon number transmitted per symbol can be calculated by

$$\frac{1}{n} \int_0^{nT} |S(t)|^2 dt \quad [\text{transmitted photons/symbol}]. \quad (1.8)$$

After attenuation through the quantum channel, the mean photon number per symbol

of the received signal becomes η fraction of the transmitted photons, i.e.,

$$\eta \cdot \left(\frac{1}{n} \int_0^{nT} |S(t)|^2 dt \right) \text{ [received photons/symbol]}. \quad (1.9)$$

We consider communication under the constraint of the transmitted mean photon number,

$$\mathbb{E} \left[\frac{1}{n} \int_0^{nT} |S(t)|^2 dt \right] \leq \mathcal{E}, \quad (1.10)$$

where the expectation is taken over the codewords. When we assume a constant average *power* over a symbol time T , the average photon number (or energy) per symbol is proportional to the symbol time T . Therefore, there is a one-to-one correspondence between the symbol time and the average photon number. A shorter symbol time, T , or equivalently a smaller average photon number per symbol, \mathcal{E} , means that we need to transmit a new symbol every short time interval T . In other words, after sending a signal with a small average photon number of \mathcal{E} , we need to immediately change the signal according the next input symbol. Since optical devices modulate the physical signal to transmit over quantum channels, there are inherent physical constraints on decreasing \mathcal{E} , or equivalently modulating the input signal in a fast way. Moreover, decreasing the input symbol period, T , or equivalently decreasing \mathcal{E} , results in a low spectral efficiency (bits/sec./Hz), since it results in wide-band communications. Therefore, it is important to note that there are limitations in decreasing \mathcal{E} , caused by optical devices as well as low spectral efficiencies.

At the receiver side, if we measure the optical wave $E_r(t)$ by a photon counter, which counts the number of photon arrivals during each symbol period of T , and denote the number of photon arrivals during the i -th symbol time as $Y_i \in \{0, 1, \dots, \infty\}$, then the output symbol distribution follows the Poisson process with rate $S_0 := \eta \cdot \left(\int_0^T |S_0(t)|^2 dt \right)$ or $S_1 := \eta \cdot \left(\int_0^T |S_1(t)|^2 dt \right)$, when the i -th input symbol is $X_i = 0$ or $X_i = 1$, respectively.

Thus, the probability distribution of Y_i , given X_i , is

$$P(Y_i = y | X_i = x) = \frac{(S_x)^y e^{-S_x}}{y!}, \quad x = 0, 1. \quad (1.11)$$

After detecting the received optical signal $S(t)$, by a photon counter, we get n -output symbols $Y^n \in \mathcal{Y}^n$, which are then decoded by $g : Y^n \rightarrow \{1, \dots, e^{nR}\}$.

We can rewrite the maps of modulator and free space optical channel using Dirac notation (Bra-ket notation). We can denote the signal wave $\{S_0(t), S_1(t)\}$, $t \in [0, T)$, corresponding to the input symbol $X \in \{0, 1\}$ as coherent states $\{|\mathbf{S}_0\rangle, |\mathbf{S}_1\rangle\}$ using ket notation, where bold face is used to emphasize the fact that input states are not general quantum states, but *coherent states* (ideal laser light). The coherent state $|\gamma\rangle$ can be thought of as a vector in a Hilbert space \mathcal{H} , and the average photon number of a coherent state $|\gamma\rangle$ is equal to $|\gamma|^2$. The modulator can be written as a map $X \rightarrow |\mathbf{S}_X\rangle$, $X \in \{0, 1\}$, and the free space optical channel as a map $|\mathbf{S}_X\rangle \rightarrow |\eta\mathbf{S}_X\rangle$. Therefore, we can define a new *effective* mapping from the input symbol $X \in \mathcal{X}$ to channel output state $|\eta\mathbf{S}_X\rangle$, which includes all the effects of modulation and optical channel, as $W : x \rightarrow |\eta\mathbf{S}_x\rangle$. The measurement of this quantum state by a photon counter can also be described by a POVM $\{\Pi_y\}$, $y \in \{0, 1, \dots, \infty\}$, with $|\Pi_y\rangle = |y\rangle\langle y|$ where $|y\rangle$ is a (photon) number state (Fock state).

Now we can ask questions as follows: What is the maximum achievable information rate (under the mean photon number constraint of \mathcal{E}) when we measure each coherent state by a photon counter? Can we achieve the capacity of this channel, which is written in (1.1), with a simple photon counter?

In [6], we show that the maximum achievable information rate with a photon counter (direct detection) is

$$C_{\text{DD}}(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}) \text{ [nats/symbol]} \quad (1.12)$$

as $\mathcal{E} \rightarrow 0$, which falls short of approaching the capacity of the bosonic channel,

$$C(\mathcal{E}) = (\mathcal{E} + 1) \log(\mathcal{E} + 1) - \mathcal{E} \log \mathcal{E} = \mathcal{E} \log \frac{1}{\mathcal{E}} + \mathcal{E} + o(\mathcal{E}) \text{ [nats/symbol]} \quad (1.13)$$

as $\mathcal{E} \rightarrow 0$. For the calculation of these capacities, it is assumed that the attenuation parameter over the communication link, $\eta = 1$. This is one example showing that the capacity of a quantum channel, which is proven to be achievable in quantum information theory, is not easily achievable with practical quantum receivers of optical communication. Even though the two capacities make a difference only in the 2^{nd} dominant term, this difference matters greatly in the practical design of optical communication systems when we consider the trade-offs between photon information efficiency (bits/photon) and spectral efficiency (bits/sec./Hz), as will be shown in Chapter 4. These discrepancies between theoretical limits and practically achievable performances are observed very often in quantum communication/detection problems.

That motivates us to ask new questions as follows: under *practical* assumptions on the quantum receiver, how well can the system perform (in terms of information rate for communication problems, or detection error probability for quantum detection problems)? What properties of *specific* sets of quantum measurements in fact generate improvements in performance, by increasing the efficiency in extracting information from quantum states? In this thesis, we will shed some light on these new kinds of questions.

■ 1.3 Adaptive Measurement

Our first consideration regards adaptive measurement, which is the most general form of the length-1 quantum measurement. With adaptive measurement, we measure each received quantum state one at a time, and then update the next measurement based on the previous observations. Compared to the joint measurements over a large number of quantum states, it is much easier to physically implement adaptive measurement.

Assume that we receive a sequence of length- n unknown quantum states, $|\psi_{x_1}\rangle \otimes |\psi_{x_2}\rangle \otimes \cdots \otimes |\psi_{x_n}\rangle \in \mathcal{H}^{\otimes n}$, for $x_i \in \mathcal{X}$, $i = 1, \dots, n$. We want to identify (x_1, x_2, \dots, x_n) by measuring these quantum states with an adaptive measurement. We measure each state one at a time, and then design the next measurement based on the outputs of the previous observations; the i -th received state, $|\psi_{x_i}\rangle$, is measured by a POVM $\{\Pi_{y_i}(y_1^{i-1})\}$, $y_i \in \mathcal{Y}_i$, that depends on the previous observations y_1^{i-1} . The resulting channel distribution can be written as

$$p(y_1^n | x_1^n) = \prod_{i=1}^n p(y_i | x_i, y_1^{i-1}) = \prod_{i=1}^n \langle \psi_{x_i} | \Pi_{y_i}(y_1^{i-1}) | \psi_{x_i} \rangle. \quad (1.14)$$

Note that the resulting channel distribution is not as general as the one in (1.3), which can be generated by the joint measurement over n -symbols; however, this channel can have *memory*, which may help increase the efficiency in extracting information from quantum states, compared to that of the memoryless channel in (1.4), which is generated by (non-adaptive) product measurements.

Then, how can we utilize this ability of adaptive measurement to increase the efficiency in extracting classical information from quantum states? What is the best performance we can achieve with adaptive measurement in quantum detection and communication problems?

There have been many attempts to show the performance bounds on adaptive measurements [1, 10, 40]. For example, in [10], the binary hypothesis testing between two coherent states (ideal laser lights) with adaptive measurement is considered. According to the two hypotheses $H = 0, 1$, either the complex waveform $S_0(t)$ or $S_1(t)$, $t \in [0, T]$, is transmitted with prior probabilities $\{p_0, p_1\}$ respectively. The theoretical lower bound of detection error probability (the Helstrom bound) of this problem over all possible quantum detectors is calculated in [21] as

$$P_e = \frac{1}{2} \left(1 - \sqrt{1 - 4p_0p_1 e^{-\int_0^T |S_0(t) - S_1(t)|^2 dt}} \right). \quad (1.15)$$

Surprisingly, this theoretical lower bound turns out to be achievable even with a very simple receiver structure (the so-called Dolinar receiver), which can adaptively update its feedback control signal based on previous photon arrivals. Therefore, for the binary hypothesis problem between two coherent states, we can achieve the theoretical lower bound of detection error probability with adaptive measurement. Unfortunately, this result does not generalize to detection problems with more than two hypotheses.

However, in [4], an example of 9-ary hypothesis testing among length-2 *orthogonal* states has been provided, in which the states can be perfectly distinguished with a length-2 joint measurement, but can never be reliably distinguished by any type of product measurements, even if the observers are allowed to update the next measurement based on the previous observations. Therefore, we know that the answer for the question of whether or not there exists an adaptive measurement that can perform as well as the optimal joint measurement varies, depending on the number of hypotheses and the possible input states for the hypothesis testing. However, there had been no general theory to show for which problems, there exist adaptive measurements that can perform as well as the optimum joint measurement, and for which, there does not.

In this thesis, we derive the necessary and sufficient conditions for adaptive measurement to perform as well as the optimum joint measurement for M -ary hypothesis testing problems. This result can be either a guide to derive the optimal adaptive measurement, or to prove the nonexistence of the adaptive measurement achieving the theoretical lower bound of detection error probability (the Helstrom bound) for a certain class of hypothesis testing. By using this result, we show that the Dolinar receiver in [10] is the exact physical translation of the mathematical description for the optimal adaptive measurement that satisfies the necessary and sufficient conditions to achieve the Helstrom bound in (1.15).

We also re-derive the Dolinar receiver for binary hypothesis testing, with the aim of maximizing communication efficiency at each instant, based on recursively updated knowledge of the receiver. Using this viewpoint, we give a natural generalization of the

Dolinar receiver to general M -ary hypothesis testing problems. We also analyze the information capacity with adaptive measurement, and compare the result with that of direct detection receivers and arbitrary quantum receivers (the Holevo limit), using the appropriate scalings in the low photon number regime.

■ 1.4 Joint Measurement: Superadditivity of Quantum Channel Coding Rate

We then turn our attention to the finite blocklength joint detection receiver (JDR). In particular, we investigate the maximum achievable information rate with finite blocklength quantum measurements. The length of quantum measurements N is now independent of the overall length of the quantum codeword, which is denoted as $N_c \geq N$. The JDR measures each sub-block of N received symbols. In order to reliably decode the message encoded over N_c symbols, the receiver collects N_c/N classical outputs generated by measuring each N -symbol sub-block, and then applies the optimum *classical* decoding algorithm over the collected outputs from many such sub-blocks.

The overall operation, depicted in Fig. 1.2, is a *concatenated coding* system. Each sub-block of N symbols is generated by an inner code of length N and rate R . There is also an outer code of length $n = N_c/N$ and rate r that is decoded by a classical outer decoder. The inner encoder, the quantum channel, and the quantum joint-detection receiver can be collectively viewed as a discrete memoryless *superchannel* with e^{NR} inputs and outputs, with transition probabilities $p_{k|j}^{(N)}$, $j, k \in \{1, \dots, e^{NR}\}$ induced by the choice of the inner code and the JDR that collectively measures sequences of N quantum symbols¹. We denote the maximum mutual information of the superchannel attainable by an optimal choice of a length N inner code and a length N JDR, as:

$$C_N := \max_{p_j} \max_{\{N\text{-symbol inner code-measurement pairs}\}} I(p_j, p_{k|j}^{(N)}) \quad (1.16)$$

¹Note that this inner joint-detection quantum decoder may in principle need more than e^{NR} classical outcomes to attain the maximum mutual information possible over all choices of inner codes and quantum measurements.

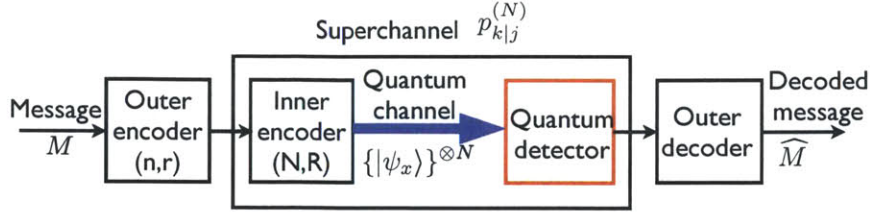


Figure 1.2. Concatenated coding over a classical-quantum channel

where $I(p_j, p_{k|j}^{(N)})$ is the classical mutual information of the channel distribution $p_{k|j}^{(N)}$ with the input distribution p_j .

As the outer code length n goes to infinity, we can reliably send e^{nr} messages as long as $r < C_N$. Since the length of the overall concatenated code is $N_c = nN$, the maximum achievable information rate at a finite N becomes $(\log e^{nr})/(nN) = r/N < C_N/N$. Therefore, C_N/N is the maximum achievable information rate at a finite length N of joint measurements. The question we pose is: How does C_N/N scale with N ?

By definition, C_N is superadditive in N , i.e., $C_N + C_M \leq C_{N+M}$, and its limit $\lim_{N \rightarrow \infty} C_N/N = C$ where C is the Holevo capacity [23]. Moreover, it is known that for some examples of input states, strict superadditivity of C_N can be demonstrated [24, 33]. For example, for a binary channel with inputs $\{|\psi_0\rangle, |\psi_1\rangle\}$ with $|\langle\psi_0|\psi_1\rangle| = \gamma$,

$$\begin{aligned}
 C &= -[(1-\gamma)/2] \log [(1-\gamma)/2] - [(1+\gamma)/2] \log [(1+\gamma)/2] \\
 C_1 &= \left[(1-\sqrt{1-\gamma^2}) \log (1-\sqrt{1-\gamma^2}) + (1+\sqrt{1-\gamma^2}) \log (1+\sqrt{1-\gamma^2}) \right] / 2.
 \end{aligned}
 \tag{1.17}$$

We can see that $C_1 < C$ for every $0 < \gamma < 1$, meaning C_N/N is strictly increasing in N , with limit equal to C as $N \rightarrow \infty$.

However, the calculation of C_N , even for the binary alphabet, is extremely hard for $N > 1$ because the complexity of optimization increases exponentially with N . In this thesis, we provide a lower bound on C_N for finite N , for quantum channels. A new framework for understanding the strict superadditivity of C_N in quantum channels will also be provided, which is different from the previous explanation of the phenomenon

by *entangling measurements* and the resulting memory in the quantum channel [35]. Moreover, under the new framework, the superadditivity of C_N , which has been mostly thought to be a unique property observed in quantum channels but not in classical discrete memoryless channels (DMCs), can be understood as a more general phenomenon that happens even in classical DMCs when the concatenated code is used with an inner decoder that makes a hard-decision at a finite inner code blocklength.

■ 1.5 Thesis Outline

The rest of this thesis is divided into four major parts. In Chapter 2, we will briefly summarize background on quantum information theory, with the focus on the quantum hypothesis testing, and communication (sending classical bits) over quantum channels. We will present the previously known theoretical bounds on performance for the quantum detection and communication problems. For the quantum detection problem, we will summarize the results in [21], which prove the theoretical lower bound (the Helstrom bound) on the average probability of detection error for M -ary hypothesis testing over quantum states. We also present the Holevo capacity for classical communication over quantum channels [20, 23], the maximum achievable information rate over quantum channels, under the assumptions of infinite uses of quantum channels as well as the ability to jointly measure the overall block of quantum states.

In Chapter 3, we study the quantum detection problems with adaptive measurements. We derive the necessary and sufficient conditions for the adaptive measurement to perform as well as the optimum joint measurement. Based on this result, we derive the adaptive measurement that achieves the Helstrom limit for the binary hypothesis between n -copy of quantum states. Moreover, we show that the Dolinar receiver, which is a simply structured receiver that has been known to achieve the Helstrom limit for discrimination between binary coherent states, exactly translates what the optimum adaptive measurement mathematically describes into a realizable receiver with only a photon counter and feedback control system.

In Chapter 4, we show that for the binary hypothesis testing between two ideal laser light pulses, if we update the adaptive measurement to maximize the communication efficiency at each instant, based on recursively updated knowledge of the receiver, then we can perform as well as in the case when we can collectively measure the received laser light of an entire duration. In this viewpoint, the Dolinar receiver for the binary hypothesis testing can be re-derived. Based on this viewpoint, we consider the generalization of the Dolinar receiver for M -ary hypothesis problems. We also analyze capacity with coherent receivers (generalized Dolinar receivers), and compare the results with those for direct detection receivers and arbitrary quantum receivers (the Holevo limit), using the appropriate scalings in the low photon number regime.

In Chapter 5, we analyze superadditivity—the phenomenon that the maximum accessible information rate per channel use strictly increases as the number of channel outputs jointly measured at the receiver increases—at the maximum information rate achievable over a pure-state classical-quantum channel. We analyze the rate vs. complexity trade-off by considering the capacity of the classical discrete memoryless superchannel induced under a concatenated coding scheme, where the quantum measurement acts exclusively on the length- N inner codewords, while allowing arbitrary classical outer-code complexity. We prove a general lower bound on the maximum accessible information per channel use for a finite-length joint measurement, and express it in terms of V , the quantum version of channel dispersion, and C , the channel capacity. The superadditivity is observed even in the capacity of a classical discrete memoryless channel (DMC) in a concatenated coding scheme due to loss of information from hard-decisions by the inner decoder over blocklength N . We develop a unifying framework, within which the superadditivity in capacity of the classical DMC and that of a classical-quantum channel can both be expressed by V/C^2 , a quantity that we show is proportional to the inner-decoder measurement length N that is sufficient to achieve a given fraction $0 < \alpha \leq 1$ of the capacity.

Chapter 6 contains our conclusions and suggestions for future work.

Background

■ 2.1 Quantum Hypothesis Testing

The theory of quantum hypothesis testing was first developed and established by Helstrom in the 1970's. In this section, we briefly summarize the important results of [21] and present a simple and concrete example of quantum binary hypothesis testing.

Imagine an optical communication system transmitting messages chosen from an alphabet of M different symbols. For each symbol in the alphabet, the transmitter produces a different optical signal of duration T seconds. The optical signals are modulations of the electromagnetic field, produced by a laser and lasting no longer than T seconds.

A receiver detects the light incident on its aperture A . We suppose that the receiver is synchronized with the transmitter so that it can identify the beginning and end of the transmitted signal. During the interval, the receiver observes the light incident and based on this observation decides which one of the M signals was transmitted.

The optical signal that arrives at the aperture of the receiver during the observation interval $[0, T)$ can be described by a set of M *density operators* $\rho_1, \rho_2, \dots, \rho_M$ for each of the M possible transmitted signals. We denote the Hilbert space where M different density operators stay by \mathcal{H} . The observation procedure at the receiver can be specified by a Positive Operator-Valued Measure (POVM), a set of M *detection operators*

$\Pi_1, \Pi_2, \dots, \Pi_M$. These operators should satisfy

$$\sum_{j=1}^M \Pi_j = \mathbb{1}; \quad \Pi_j \geq 0 \quad \text{for } \forall j \in \{1, \dots, M\}, \quad (2.1)$$

i.e., they are non-negative Hermitian operators resolving identity in \mathcal{H} .

Let us consider the M -ary hypothesis testing problem. There are M hypotheses about the state of a quantum system, where the i th hypothesis H_i indicates that the system is in the state ρ_i . When the quantum system is observed by a set of detection operators $\Pi_1, \Pi_2, \dots, \Pi_M$, one and only one of the outputs is going to click. When we observe a click for a detection operator Π_j , our guess for the true hypothesis becomes H_j . From the laws of quantum mechanics, the conditional probabilities that the output corresponding to Π_j clicks when H_i is true is

$$\Pr(j|i) = \text{Tr}(\Pi_j \rho_i), \quad i, j = 1, \dots, M. \quad (2.2)$$

Let us denote the prior probabilities of hypothesis H_i by p_i for $i = 1, \dots, M$. The cost of choosing hypothesis H_j when H_i is true is denoted as C_{ji} . Then the average cost of the M -ary hypothesis testing specified by the detection operator $\{\Pi_j\}$ is

$$\bar{C} = \sum_{i,j=1}^M p_i C_{ji} \text{Tr}(\Pi_j \rho_i) = \text{Tr} \left(\sum_{j=1}^M \Pi_j W_j \right) \quad (2.3)$$

where the Hermitian risk operators W_j are defined by

$$W_j = \sum_{i=1}^M p_i C_{ji} \rho_i. \quad (2.4)$$

We want to minimize (2.3) under the constraints (2.1). This optimization resembles linear programming, except that operators are involved instead of functions. The necessary and sufficient conditions for the optimum $\{\Pi_j\}$ that minimizes (2.3) was provided in [51] as follows.

Lemma 2.1. *The optimum detection operator $\{\Pi_j\}$ should satisfy*

$$(W_j - \Gamma)\Pi_j = \Pi_j(W_j - \Gamma) = 0, \quad j = 1, 2, \dots, M, \quad (2.5)$$

$$W_j - \Gamma \geq 0, \quad j = 1, 2, \dots, M, \quad (2.6)$$

with the Lagrange operator

$$\Gamma = \sum_{i=1}^M \Pi_i W_i = \sum_{i=1}^M W_i \Pi_i \quad (2.7)$$

required to be Hermitian.

For a special cost function of $C_{ij} = -\delta_{ij}$ where

$$\delta_{ij} := \begin{cases} 1, & i = j; \\ 0, & i \neq j, \end{cases} \quad (2.8)$$

the average cost function \bar{C} is equal to $-1 + p_e$ where p_e is the average probability of detection error, i.e.,

$$p_e = 1 - \sum_{i=1}^M p_i \text{Tr}(\Pi_i \rho_i). \quad (2.9)$$

Therefore, for this cost function, a set of detection operators that minimizes \bar{C} does minimize the average probability of error, p_e . The equations for the optimum detection operators (2.5)–(2.7) can be simplified for this cost function as

$$\Pi_j(\Gamma - p_j \rho_j) = 0, \quad j = 1, 2, \dots, M, \quad (2.10)$$

$$\Gamma - p_j \rho_j \geq 0, \quad j = 1, 2, \dots, M, \quad (2.11)$$

$$\Gamma = \sum_{i=1}^M p_i \Pi_i \rho_i = \sum_{i=1}^M p_i \rho_i \Pi_i, \quad (2.12)$$

with the change of sign of the Lagrange operator Γ .

For general quantum states, it is very difficult to solve the equations for the optimum

detection operator (2.5)–(2.7). In the rest of this section, we focus on the case where the quantum system is in a pure state, i.e., $|\psi_i\rangle \in \mathcal{H}$ under each hypothesis H_i , $i = 1, 2, \dots, M$, and provide important properties of the optimum detection operators for this case. The density operators for pure states are

$$\rho_i = |\psi_i\rangle\langle\psi_i|. \quad (2.13)$$

When these M vectors $|\psi_i\rangle$ are linearly independent, they span an M -dimensional subspace \mathcal{H}_M of the Hilbert space \mathcal{H} of the quantum system. Even though the optimum detection operator $\{\Pi_j\}$ should be a set of non-negative Hermitian operators resolving the identity in \mathcal{H} , these operators can also be confined to operating on \mathcal{H}_M since the components of Π_i outside of \mathcal{H}_M do not change the conditional probability $\Pr(j|i)$.

To verify this argument, let G_M be a projector onto \mathcal{H}_M , and write each operator of $\{\Pi_j\}$ as

$$\Pi_j = G_M \Pi_j G_M + (1 - G_M) \Pi_j G_M + G_M \Pi_j (1 - G_M) + (1 - G_M) \Pi_j (1 - G_M). \quad (2.14)$$

Then

$$\Pr(j|i) = \text{Tr}(\rho_i \Pi_j) = \langle\psi_i|\Pi_j|\psi_i\rangle = \text{Tr}(\rho_i \Pi'_j), \quad \text{where } \Pi'_j := G_M \Pi_j G_M \quad (2.15)$$

since $(1 - G_M)|\psi_k\rangle = \langle\psi_k|(1 - G_M) = 0$.

The condition for the detection operator to resolve identity in \mathcal{H} can thus be generalized to

$$\sum_{j=1}^M \Pi'_j + (1 - G_M) = \mathbb{I}, \quad (2.16)$$

for $\{\Pi'_j\}$ in \mathcal{H}_M . Since $(1 - G_M)$ has no effect on the conditional property for the detection of the pure states in \mathcal{H}_M , we can confine our analysis to the subspace \mathcal{H}_M spanned by the pure states $\{|\psi_i\rangle\}$.

For M linearly independent pure states $\rho_i = |\psi_i\rangle\langle\psi_i|$, when the simple cost function

$C_{ij} = -\delta_{ij}$, which minimizes the average error probability, is adopted, the optimum detection operator $\{\Pi_i\} \in \mathcal{H}_M$ satisfying the conditions (2.10)–(2.12) has a simple structure of orthonormal projectors in \mathcal{H}_M . This fact was first proved by Kennedy in [27]. Moreover, it was shown that such orthonormal projectors are uniquely determined in \mathcal{H}_M [28].

Lemma 2.2. *For M linearly independent pure states $\rho_i = |\psi_i\rangle\langle\psi_i|$, the optimum detection operator $\{\Pi_j\}$ that minimizes the average error probability is the set of M orthonormal projectors, i.e.,*

$$\Pi_j = |\omega_j\rangle\langle\omega_j| \quad (2.17)$$

where $\{|\omega_j\rangle\}$ is a set of orthonormal vectors spanning \mathcal{H}_M . Therefore, it satisfies

$$\Pi_i\Pi_j = \delta_{ij}\Pi_i = \delta_{ij}\Pi_j. \quad (2.18)$$

The set of orthonormal measurement vectors $\{|\omega_j\rangle\}$ satisfying (2.10)–(2.12) is uniquely determined.

Even though this result does not directly provide a solution $\{\Pi_j\}$ for the equations (2.10)–(2.12), this makes it much easier to find the optimum detection operator since we can now assume that $\{\Pi_j\}$ are orthonormal projectors.

Based on Lemma 2.2, it is straightforward to find the optimum measurement vectors that minimize the probability of error for binary hypothesis testing problems. When $\rho_0 = |\psi_0\rangle\langle\psi_0|$ and $\rho_1 = |\psi_1\rangle\langle\psi_1|$ for hypotheses H_0 and H_1 , respectively, the pure states can be written as

$$\begin{aligned} |\psi_0\rangle &= \cos\theta|x\rangle + \sin\theta|y\rangle, \\ |\psi_1\rangle &= \cos\theta|x\rangle - \sin\theta|y\rangle. \end{aligned} \quad (2.19)$$

for the orthonormal basis $\{|x\rangle, |y\rangle\}$ of the 2-dimensional subspace spanned by the two input states $\{|\psi_0\rangle, |\psi_1\rangle\}$.

Then, the orthonormal projectors $\Pi_0 = |\omega_0\rangle\langle\omega_0|$ and $\Pi_1 = |\omega_1\rangle\langle\omega_1|$ can be parame-

terized by ϕ and can be expressed as

$$\begin{aligned} |\omega_0\rangle &= \cos \phi |x\rangle + \sin \phi |y\rangle, \\ |\omega_1\rangle &= \cos\left(\phi - \frac{\pi}{2}\right) |x\rangle + \sin\left(\phi - \frac{\pi}{2}\right) |y\rangle = \sin \phi |x\rangle - \cos \phi |y\rangle. \end{aligned} \quad (2.20)$$

The average probability of error for these measurement vectors is

$$\begin{aligned} p_e &= p_0 |\langle \psi_0 | \omega_1 \rangle|^2 + p_1 |\langle \psi_1 | \omega_0 \rangle|^2 \\ &= p_0 (\cos \theta \sin \phi - \sin \theta \cos \phi)^2 + p_1 (\cos \theta \cos \phi - \sin \theta \sin \phi)^2 \\ &= p_0 \sin^2(\phi - \theta) + p_1 \cos^2(\phi + \theta). \end{aligned} \quad (2.21)$$

By taking the partial derivative of p_e with respect to ϕ , it can be shown that the optimum ϕ^* should satisfy

$$\frac{p_1}{p_0} = \frac{\cos(\phi^* - \theta) \sin(\phi^* - \theta)}{\cos(\phi^* + \theta) \sin(\phi^* + \theta)} = \frac{\sin 2\phi^* \cos 2\theta - \sin 2\theta \cos 2\phi^*}{\sin 2\phi^* \cos 2\theta + \sin 2\theta \cos 2\phi^*}, \quad (2.22)$$

which is equivalent to

$$\begin{aligned} \sin 2\phi^* &= \frac{(p_0 + p_1) \sin 2\theta}{\sqrt{(p_0 + p_1)^2 \sin^2 2\theta + (p_0 - p_1)^2 \cos^2 2\theta}} = \frac{\sin 2\theta}{\sqrt{1 - 4p_0 p_1 \cos^2 2\theta}}, \\ \cos 2\phi^* &= \frac{(p_0 - p_1) \sin 2\theta}{\sqrt{(p_0 + p_1)^2 \sin^2 2\theta + (p_0 - p_1)^2 \cos^2 2\theta}} = \frac{(p_0 - p_1) \cos 2\theta}{\sqrt{1 - 4p_0 p_1 \cos^2 2\theta}}. \end{aligned} \quad (2.23)$$

The minimum average error probability at ϕ^* is

$$p_e^* = \frac{1 - \sqrt{1 - 4p_0 p_1 \cos^2 2\theta}}{2} = \frac{1 - \sqrt{1 - 4p_0 p_1 |\langle \psi_0 | \psi_1 \rangle|^2}}{2}. \quad (2.24)$$

For M -ary hypothesis testing, $M > 2$, even for pure states, it is very difficult to find the explicit form for the optimum measurement vectors and to calculate the minimum average error probability. Only for the special cases where there exists a certain symmetry between the pure states, was it known how to find the optimum measurements and the resulting p_e [21].

■ 2.2 Capacity of Classical-Quantum Channel: Holevo capacity

In this section, we consider the problem of sending classical bits over quantum channels $W : x \rightarrow |\psi_x\rangle$, $x \in \mathcal{X}$ where $\{|\psi_x\rangle\} \in \mathcal{H}$. We will first present the quantum version of typical subspaces and of the asymptotic equipartition property (AEP). Based on these tools, we will present the achievability proof of the Holevo capacity, first shown in [20, 23].

Suppose we use a quantum channel $W : x \rightarrow |\psi_x\rangle$ with $x \in \mathcal{X}$ chosen with probabilities p_x . The density operator of the output of this quantum channel becomes $\rho = \sum_{x \in \mathcal{X}} p_x |\psi_x\rangle\langle\psi_x|$ in \mathcal{H} . For N uses of this quantum channel, the density operator for the N -sequence of channel outputs can be written as $\rho^{\otimes N} = \rho \otimes \dots \otimes \rho$ in $\mathcal{H}^{\otimes N}$. The N -fold tensor product Hilbert space $\mathcal{H}^{\otimes N}$ can be decomposed into two subspaces: a “typical” subspace Λ and the perpendicular subspace Λ^\perp . The typical subspace is defined as a subspace spanned by a set of eigenstates of $\rho^{\otimes N}$ whose eigenvalues γ_i satisfy

$$2^{-N(H(\rho)+\delta)} < \gamma_i < 2^{-N(H(\rho)-\delta)} \quad (2.25)$$

for an arbitrarily small $\delta > 0$, when $H(\rho)$ is the von Neumann entropy of the density operator ρ , i.e., $H(\rho) = -\text{Tr}(\rho \log \rho)$. Note that this definition of typical subspace resembles that of classical typical subspace, when we interpret the eigenvalues of $\rho^{\otimes N}$ as the probabilities of observing the output sequence corresponding to each eigenstate. From this definition of typical subspace, we can derive the following three properties. Let Π_Λ denote the projector onto the typical subspace Λ , and $\dim\Lambda$ denote the number of dimensions in the typical subspace.

1. For an eigenstate in the typical subspace Λ , its eigenvalue γ satisfies

$$H(\rho) - \delta < -\frac{1}{N} \log \gamma < H(\rho) + \delta.$$

2. For an arbitrarily small $\epsilon > 0$, $\text{Tr}(\Pi_\Lambda \rho^{\otimes N} \Pi_\Lambda) > 1 - \epsilon$, for N sufficiently large.

3. The dimension of Λ satisfies

$$(1 - \epsilon)2^{N(H(\rho) - \delta)} \leq \dim \Lambda \leq 2^{N(H(\rho) + \delta)}. \quad (2.26)$$

The first property can be directly derived from the definition of the typical subspace in (2.25). The second property can be proven by using the law of large numbers. The third property can be derived from the first and the second properties. We will use these properties to prove the achievability of the Holevo capacity.

The following theorem summarizes the result of the classical capacity of quantum channels $W : x \rightarrow |\psi_x\rangle$, shown in [20, 23].

Theorem 2.3. *Consider a classical-quantum channel $W : x \rightarrow |\psi_x\rangle$. For every rate $R < C$, there exists a length- N and rate- R code that can be decoded by a set of length- N joint measurements with $P_e \rightarrow 0$ as $N \rightarrow \infty$. The converse has also been proved in [22]. The Holevo capacity C is*

$$C = \max_{p_x} -\text{Tr}(\rho \log \rho) = \max_{p_x} H(\rho)$$

where $\rho = \sum_{x \in \mathcal{X}} p_x |\psi_x\rangle\langle\psi_x|$.

Proof. We present the proof of this theorem from [20]. Let us first introduce some notations related to the quantum codewords of length N and rate R codes. First, denote the total number of messages as $M := 2^{NR}$. The encoder $f : \{1, \dots, M\} \rightarrow \mathcal{X}^N$ maps each message into a length- N codeword. The codeword for the j -th message can be written as $f(j) = (x_1(j), \dots, x_N(j))$ where $x_i(j) \in \mathcal{X}$ for $i = 1, \dots, N$. After sending each coded symbol through the classical-quantum channel $W : x \rightarrow |\psi_x\rangle$, the received length- N sequence of states can be written in a density operator form,

$$S_{f(j)} := |\psi_{x_1(j)}\rangle\langle\psi_{x_1(j)}| \otimes \cdots \otimes |\psi_{x_N(j)}\rangle\langle\psi_{x_N(j)}|. \quad (2.27)$$

When we denote $|\psi_{f(j)}\rangle := |\psi_{x_1(j)}\rangle \otimes \cdots \otimes |\psi_{x_N(j)}\rangle$,

$$S_{f(j)} = |\psi_{f(j)}\rangle\langle\psi_{f(j)}|. \quad (2.28)$$

Let us define a matrix Ψ such that its j -th column is $|\psi_{f(j)}\rangle$. When $|\psi_{f(j)}\rangle$, $j \in \{1, \dots, M\}$ stay in a d -dimensional Hilbert space, the singular value decomposition of Ψ can be represented as

$$\Psi = (|\psi_{f(1)}\rangle, |\psi_{f(2)}\rangle, \dots, |\psi_{f(M)}\rangle) = U\Sigma V^\dagger, \quad (2.29)$$

where U and V are unitary matrices of size $d \times d$ and $M \times M$, respectively, and Σ is a $d \times M$ rectangular diagonal matrix with non-negative real numbers on the diagonal. The operator V^\dagger is the Hermitian conjugate of V . We also denote the Gram matrix, Γ , and the Gram operator, G , of Ψ as

$$\begin{aligned} \Gamma &= \Psi\Psi^\dagger = V(\Sigma^\dagger\Sigma)V^\dagger, \\ G &= \Psi^\dagger\Psi = U(\Sigma\Sigma^\dagger)U^\dagger. \end{aligned} \quad (2.30)$$

Note that Γ and G are positive operators.

Now, we will introduce Square Root Measurements (SRM) $\{\Pi_j\}$, with which the M encoded quantum states are measured. The SRM is defined as a rank-one operator such that

$$\Pi_j = |\omega_j\rangle\langle\omega_j| = \left(G^{-1/2}\right) S_{f(j)} \left(G^{-1/2}\right) = \left(G^{-1/2}|\psi_{f(j)}\rangle\right) \left(\langle\psi_{f(j)}|G^{-1/2}\right) \quad (2.31)$$

where

$$G^{-1/2} = U \left((\Sigma\Sigma^\dagger)^{-1/2} \right) U^\dagger, \quad (2.32)$$

and $(\Sigma\Sigma^\dagger)^{-1/2}$ is formed by replacing every non-zero diagonal entry of $\Sigma\Sigma^\dagger$ with one over the square root of each entry. Note that the defined measurement $\{\Pi_j\}$ satisfies $\Pi_j \geq 0$, for all $j \in \{1, \dots, M\}$, and $\sum_{j=1}^M \Pi_j = \mathbb{1}$.

When we define a matrix Ω such that its j -th column is the j -th measurement vector $|\omega_j\rangle$,

$$\Omega = (|\omega_1\rangle, |\omega_2\rangle, \dots, |\omega_M\rangle) = G^{-1/2}\Psi = U(\Sigma\Sigma^\dagger)^{-1/2}\Sigma V^\dagger = U\Sigma(\Sigma^\dagger\Sigma)^{-1/2}V^\dagger. \quad (2.33)$$

Then, (k, j) -entry of $\Omega^\dagger\Psi$ becomes $\langle\omega_k|\psi_{f(j)}\rangle$ and

$$\Omega^\dagger\Psi = V(\Sigma^\dagger\Sigma)^{-1/2}(\Sigma^\dagger\Sigma)V^\dagger = V(\Sigma^\dagger\Sigma)^{1/2}V^\dagger = \Gamma^{1/2}. \quad (2.34)$$

Therefore, the probability that the decoder chooses the k -th message by the SRM, when the j -th message is the true one, which is denoted as $p_{k|j}^{(N)}$, is

$$p_{k|j}^{(N)} = |\langle\omega_k|\psi_{f(j)}\rangle|^2 = |\sqrt{\Gamma}_{k,j}|^2 \quad (2.35)$$

where $\sqrt{\Gamma}_{k,j}$ denotes the (k, j) -entry of $\Gamma^{1/2}$. It means that under the SRM, once we have the geometric structure of the encoded quantum states, which is represented by its Gram matrix Γ , the distribution of the measurement outputs, given the true message, can be directly calculated from $\Gamma^{1/2}$.

From (2.35), the average probability of decoding error becomes,

$$P_e = \frac{1}{M} \sum_{j=1}^M (1 - |\sqrt{\Gamma}_{j,j}|^2) = \frac{1}{M} \sum_{j=1}^M (1 - \sqrt{\Gamma}_{j,j})(1 + \sqrt{\Gamma}_{j,j}) \quad (2.36)$$

Since $\sqrt{\Gamma}_{j,j} = \langle\omega_j|\psi_{f(j)}\rangle = \langle\psi_{f(j)}|G^{-1/2}|\psi_{f(j)}\rangle$ with the positive operator $G^{-1/2}$, and $\sum_{k=1}^M |\sqrt{\Gamma}_{k,j}|^2 = 1$, the resulting $\sqrt{\Gamma}_{j,j}$, $j = 1, \dots, M$, are positive numbers in $[0, 1]$. Therefore, we can bound P_e as

$$P_e \leq \frac{2}{M} \sum_{j=1}^M (1 - \sqrt{\Gamma}_{j,j}) = \frac{2}{M} \left(M - \text{Tr} \left(\Gamma^{1/2} \right) \right). \quad (2.37)$$

When λ_i denotes the eigenvalues of $\Gamma = \Psi\Psi^\dagger$ with the decreasing order, i.e., $\lambda_1 \geq \lambda_2 \geq$

$\dots \geq \lambda_M,$

$$P_e \leq \frac{2}{M} \left(M - \text{Tr} \left(\Gamma^{1/2} \right) \right) = \frac{2}{M} \sum_{i=1}^M (\lambda_i - \sqrt{\lambda_i}), \quad (2.38)$$

since $\sum_{i=1}^M \lambda_i = M$. By using

$$\sqrt{x} \geq \frac{3}{2}x - \frac{1}{2}x^2, \quad \text{for } x \geq 0, \quad (2.39)$$

we can further bound P_e as

$$\begin{aligned} P_e &\leq \frac{2}{M} \sum_{i=1}^M \left(\lambda_i - \frac{3}{2}\lambda_i + \frac{1}{2}\lambda_i^2 \right) = \frac{1}{M} \sum_{i=1}^M (\lambda_i^2 - \lambda_i) = \frac{1}{M} \sum_{i=1}^M (\lambda_i - 1)^2 \\ &= \frac{1}{M} \sum_i \sum_{j \neq i} |\Gamma_{i,j}|^2. \end{aligned} \quad (2.40)$$

The last equality in (2.40) is from the fact that when eigenvalues of Γ are λ_i , $i = 1, \dots, M$, the eigenvalues of $\Gamma - I$ are $\lambda_i - 1$. Then, the Hilbert-Schmidt norm of $\Gamma - I$, which is $\sum_i \sum_{j \neq i} |\Gamma_{i,j}|^2$, is equal to the sum of squared eigenvalues of $\Gamma - I$.

Now, we use the random coding technique and the quantum AEP to show the existence of a length- N and rate- $R < C$ code for which the right-hand side of (2.40) converges to 0 as $N \rightarrow \infty$.

When we generate M -codewords independently according to the distribution $p_{\underline{x}} = \prod_{i=1}^N p_{x_i}$, the magnitude squared of $\Gamma_{i,j}$ averaged over the random code C is

$$\begin{aligned} \mathbb{E}_C [|\Gamma_{i,j}|^2] &= \sum_{i,j} p_{f(i)} \cdot p_{f(j)} \cdot (\langle \psi_{f(i)} | \psi_{f(j)} \rangle \langle \psi_{f(j)} | \psi_{f(i)} \rangle) \\ &= \text{Tr} \left(\sum_i p_{f(i)} |\psi_{f(i)}\rangle \langle \psi_{f(i)}| \sum_j p_{f(j)} |\psi_{f(j)}\rangle \langle \psi_{f(j)}| \right) = \text{Tr} ((\rho^{\otimes N})^2) \end{aligned} \quad (2.41)$$

where $\rho = \sum_x p_x |\psi_x\rangle \langle \psi_x|$.

Assume that the received codeword is first projected onto the typical subspace Λ of $\rho^{\otimes N}$, and then measured by the SRM defined for the projected codewords, i.e.,

$\{\Pi_\Lambda|\psi_{f(i)}\rangle\}$, $i = 1, \dots, M$. For the Gram matrix $\Gamma = \Psi\Psi^\dagger$ whose i -th column of Ψ is $\Pi_\Lambda|\psi_{f(i)}\rangle$, we can show that

$$\begin{aligned} \mathbb{E}_C[|\Gamma_{i,j}|^2] &= \sum_{i,j} p_{f(i)} \cdot p_{f(j)} \cdot (\langle\psi_{f(i)}|\Pi_\Lambda|\psi_{f(j)}\rangle\langle\psi_{f(j)}|\Pi_\Lambda|\psi_{f(i)}\rangle) \\ &= \text{Tr}(\Pi_\Lambda(\rho^{\otimes N})^2\Pi_\Lambda) < 2^{-2N(H(\rho)-\delta)} \cdot 2^{N(H(\rho)+\delta)} = 2^{-N(H(\rho)-3\delta)}, \end{aligned} \quad (2.42)$$

by using the fact that Π_Λ commutes with $\rho^{\otimes N}$, and the bound for eigenvalues in the typical space (2.25) and the dimensionality of the typical space (2.26).

Using (2.42) and (2.40),

$$\mathbb{E}_C[P_e] \leq \frac{1}{M} \sum_i \sum_{j \neq i} \mathbb{E}_C[|\Gamma_{i,j}|^2] < (M-1)2^{-N(H(\rho)-3\delta)}. \quad (2.43)$$

This bound shows that for $M = 2^{NR}$, there exists a code of length N and rate $R > H(\rho) - 3\delta$ such that it can be decoded by the SRM with P_e approaching to 0 as $N \rightarrow \infty$. Since this result is true for any input distribution p_x , it shows that the maximum achievable rate R is greater than $C - 3\delta$, i.e., $R > \max_{p_x} H(\rho) - 3\delta = C - 3\delta$ when the code blocklength $N \rightarrow \infty$. \square

Quantum Hypothesis Testing with Adaptive Measurements

■ 3.1 Introduction

Discrimination between non-orthogonal quantum states is one of the key tasks in quantum information theory for communication over quantum channels. When we have a quantum system set to be in one of the quantum states $\{|\psi_x\rangle\} \in \mathcal{H}$ with probability p_x , $x \in \mathcal{X}$, what is the optimal way to measure the system in order to generate a correct guess of the true state with minimum probability of error? And what is the resulting minimum probability of error? In [21], Helstrom answered these fundamental questions by deriving the necessary and sufficient conditions for the optimum measurements to satisfy, and he also calculated the resulting minimum probability of detection error in terms of the inner product between the possible states $\{|\psi_x\rangle\} \in \mathcal{H}$ and their probabilities p_x , $x \in \mathcal{X}$. We call this minimum error probability the Helstrom limit. Even though the necessary and sufficient conditions derived by Helstrom often do not result in a constructive answer for the optimum measurements, for the special cases such as binary hypothesis testing (BHT) or ternary hypothesis testing with symmetric states, we can indeed derive the optimum measurement vectors represented in terms of the possible set of quantum states $\{|\psi_x\rangle\}$ in the Hilbert space \mathcal{H} .

However, even for the case when we can find the optimum measurement vectors from the necessary and sufficient conditions, there still remain difficulties in implementing

a quantum receiver that can perform as mathematically described by the optimum measurements. In particular, when we consider the binary hypothesis testing with N -copy of a quantum state that is either $|\psi_0\rangle$ or $|\psi_1\rangle$, i.e., BHT between $|\psi_0^N\rangle := |\psi_0\rangle \otimes \cdots \otimes |\psi_0\rangle$ and $|\psi_1^N\rangle := |\psi_1\rangle \otimes \cdots \otimes |\psi_1\rangle$, the optimum measurement can be represented as a superposition of these N -tensor product states, and thus the optimum measurements are not product states, but *entangling measurements* (joint measurements over N -received quantum states). Since the physical implementation of entangling measurements over N quantum states is in general much complicated, the mathematical description for the optimum measurement does not directly translate into a physical realization of a receiver achieving the Helstrom limit.

As a way to resolve such difficulties in the practical implementation of quantum receivers, *adaptive measurement* has been widely considered. Adaptive measurement is a special type of product measurement with which we measure each state one at a time and use the outputs of the current and previous observations to update the next measurement. Since it is much easier to implement product measurements compared to entangling measurements, if any kinds of product measurements, including adaptive measurement, can perform as well as the optimal entangling measurement derived by Helstrom's conditions, then we prefer to use the adaptive measurement to implement quantum receivers. The question is then whether there exists an adaptive measurement that can achieve the Helstrom limit for the general M -ary hypothesis testing. This question has been partly answered for special cases. For example, it has been known that it is indeed possible to achieve the Helstrom limit with an adaptive measurement for the previously introduced BHT problem for the N -copy of a quantum state [1]. Acin *et al.* in [1] showed that an adaptive measurement that minimizes the interim probability of error step by step achieves global optimality for the BHT, i.e., the Helstrom limit. However, in [4], an example of 9-ary hypothesis testing between a set of length-2 *orthogonal* states is provided, which can be perfectly distinguished with a joint measurement over the length-2 states, but can never be reliably distinguished by any type of product

measurements, even if the observers are allowed to update the next measurement based on the previous observation, i.e., with adaptive (product) measurement. Therefore, we know that the answer for the question of whether or not there exists an adaptive measurement that can perform as well as the optimal entangling measurement varies depending on the number of hypotheses and the possible input states for the hypothesis testing. However, there had been no general theory to show for which problems, the adaptive measurement can achieve the Helstrom limit, for which, it cannot.

In this chapter, we derive the necessary and sufficient conditions for the adaptive measurement to achieve the Helstrom limit. This result can be either a guide to derive the optimal adaptive measurement, or to prove the non-existence of the adaptive measurement achieving the Helstrom limit for a certain class of hypothesis testing. By using this result, we show that the adaptive measurement suggested in [1] for the binary hypothesis testing between N -copy of a quantum state indeed meets the necessary and sufficient conditions for the optimal adaptive measurement. Moreover, for general M -ary hypothesis testing, we provide an important property for the optimal adaptive measurement to satisfy in order to achieve the Helstrom limit. We connect our analysis for the optimum adaptive measurement with the well-known Dolinar receiver [10], which is shown to achieve the Helstrom limit for any two *coherent states*, by showing that the operation by the Dolinar receiver is the exact physical realization of the mathematical description for the optimal adaptive measurement.

■ 3.2 Preliminaries

In this section, we will consider M -ary hypothesis testing and state the necessary and sufficient conditions for the adaptive measurement to achieve the theoretical lower bound on the detection error probability. With prior probabilities p_i , $i = 1, \dots, M$, a quantum system is set to be in state $\rho_i = |\psi_i\rangle\langle\psi_i|$ according to hypothesis H_i . The quantum states $|\psi_i\rangle$ are vectors in a Hilbert space \mathcal{H} . By measuring the quantum system, we want to discriminate between the possible input states and make a correct

guess for the true hypothesis with minimum probability of error. The quantum measurement can be written as a set of positive operators, $\{\Pi_j\}$, $j = 1, \dots, M$, that resolve identity in \mathcal{H} , i.e., $\sum_{j=1}^M \Pi_j = \mathbb{1}$. When we observe an output of the measurement Π_j , we choose H_j as a guess for the true hypothesis. We will denote such events as \widehat{H}_j , $j \in \{1, \dots, M\}$. The conditional probability of observing the measurement outcome corresponding to Π_j , when H_i is the true hypothesis, is equal to

$$p(j|i) := \Pr(\widehat{H}_j|H_i) = \text{Tr}(\Pi_j \rho_i) = \langle \psi_i | \Pi_j | \psi_i \rangle. \quad (3.1)$$

The resulting probability of detection error is then

$$p_e = \sum_i p_i (1 - p(i|i)) = 1 - \sum_{i=1}^M p_i \text{Tr}(\Pi_i \rho_i). \quad (3.2)$$

We need to find a set of positive operators $\{\Pi_j\}$ that minimizes p_e under the constraint of $\Pi_j \geq 0$, $j \in \{1, \dots, M\}$, and $\sum_{j=1}^M \Pi_j = \mathbb{1}$. In [21], a set of equations for the solution of this optimization is provided as:

$$\Pi_j(\Gamma - p_j \rho_j) = 0, \quad j = 1, 2, \dots, M, \quad (3.3)$$

$$\Gamma - p_j \rho_j \geq 0, \quad j = 1, 2, \dots, M, \quad (3.4)$$

for $\Gamma = \sum_{i=1}^M p_i \Pi_i \rho_i = \sum_{i=1}^M p_i \rho_i \Pi_i$.

Note that $\{\Pi_j\}$ that satisfies these conditions may not be unique, and it is not straightforward to derive a closed-form solution for the optimum measurement operator $\{\Pi_j\}$ from these conditions.

However, for the case when the state vectors $\{|\psi_i\rangle\} \in \mathcal{H}$ are linearly independent, a particular structure of $\{\Pi_j\}$, which makes it easier to find a closed solution satisfying (3.3)–(3.4), has been found in [27] and [28]. When pure states $|\psi_1\rangle, |\psi_2\rangle, \dots, |\psi_M\rangle$ are linearly independent, regardless of the dimension of \mathcal{H} , these states stay in the M -dimensional subspace \mathcal{U}_M spanned by $\{|\psi_i\rangle\}$. For this case, the measurement operator

$\{\Pi_j\}$ can be confined to operating on \mathcal{U}_M as the components of Π_j outside of \mathcal{U}_M do not affect the conditional probabilities $p(j|i)$; when $\mathcal{P}_{\mathcal{U}_M}$ is a projector onto the subspace \mathcal{U}_M , for $\Pi'_j = \mathcal{P}_{\mathcal{U}_M}\Pi_j\mathcal{P}_{\mathcal{U}_M}$,

$$p(j|i) = \text{Tr}(\Pi_j\rho_i) = \text{Tr}(\Pi'_j\rho_i). \quad (3.5)$$

From [27], it was known that for the detection operator $\{\Pi'_j\}$, which is confined to operating on \mathcal{U}_M , to satisfy the optimality conditions (3.3)–(3.4), it should be an orthonormal projector. Moreover, such orthonormal projective measurement $\{\Pi'_j\}$ is *uniquely* determined [28]. The following lemma summarizes these important properties of the optimum $\{\Pi'_j\}$.

Lemma 3.1. *To discriminate M linearly independent pure states $|\psi_i\rangle$ with minimum error probability, the optimum detection operator $\{\Pi'_j\}$, which is confined to operate on \mathcal{U}_M , should satisfy*

$$\Pi'_i\Pi'_j = \delta_{ij}\Pi'_i, \quad (3.6)$$

and such an optimum detection operator is uniquely determined on \mathcal{U}_M .

■ 3.3 Multiple-Copy States Discrimination with Adaptive Measurements

Let us consider M -ary hypothesis testing with N -copy of a quantum state, which is set to be in $|\psi_i\rangle \in \mathcal{H}$ according to the hypothesis H_i with prior probability p_i , $i = 1, \dots, M$. Since N -copy of a quantum state $|\psi_i\rangle$ can be viewed as an N -fold tensor product state $|\psi_i^N\rangle := |\psi_i\rangle \otimes \dots \otimes |\psi_i\rangle$, this problem is equivalent to the M -ary hypothesis testing with possible states $|\psi_i^N\rangle \in \mathcal{H}^{\otimes N}$, $i = 1, \dots, M$. The optimum measurement $\{\Pi_i\}$ for this M -ary hypothesis testing should satisfy (3.3)–(3.4) for $\rho_i = |\psi_i^N\rangle\langle\psi_i^N|$. When we confine the measurement operator to the M -dimensional subspace $\mathcal{U}_M \subset \mathcal{H}^{\otimes N}$, spanned by $\{|\psi_i^N\rangle\}$, from Lemma 3.1, the optimum measurement becomes a set of orthonormal projectors, which are uniquely determined. Since the optimum measurement, which often results in entangling measurement (joint measurement), is hard to implement,

we instead consider product measurements in $\mathcal{H}^{\otimes N}$, which can be adaptively updated based on the previous observations.

Let us focus on the following scenario: we observe one state at a time with a set of M -measurement vectors in \mathcal{H} , and update the next set of measurement vectors, which will be applied to the next state, based on the output of the previous observations. When the measurement vectors applied to the n -th copy of a state are denoted as $\{|\omega_n(y_1^{n-1}1)\rangle, \dots, |\omega_n(y_1^{n-1}M)\rangle\}$ for the previous observations $y_1^{n-1} \in \{1, \dots, M\}^{n-1}$, the resulting length- N adaptive (product) measurement can be written as

$$|\Omega(y_1^N)\rangle := |\omega_1(y_1)\rangle \otimes \dots \otimes |\omega_N(y_1^N)\rangle \quad (3.7)$$

for the sequence of output observations, $y_1^N \in \{1, \dots, M\}^N$. Note that there are total M^N -possible output sequences from this adaptive measurement for the M -ary hypothesis testing. To make a decision for the M -ary hypothesis testing, we need to group those M^N -possible outputs into M different decision sets; we denote the set of y_1^N 's that results in the decision of the i -th hypothesis as D_i for $i = 1, \dots, M$. Then the resulting probability of error for such decision sets D_i is

$$p_e = 1 - \sum_{i=1}^M p_i \langle \psi_i^N | \left(\sum_{y_1^N \in D_i} |\Omega(y_1^N)\rangle \langle \Omega(y_1^N)| \right) | \psi_i^N \rangle. \quad (3.8)$$

Therefore, the measurement operator $\{\Pi_i\}$, $i = 1, \dots, M$, of the adaptive measurements with grouping of the outputs can be written as

$$\Pi_i = \sum_{y_1^N \in D_i} |\Omega(y_1^N)\rangle \langle \Omega(y_1^N)|. \quad (3.9)$$

The questions we want to answer are as follows: to achieve the Helstrom limit with the adaptive measurement in (3.9), what are the necessary and sufficient conditions that the adaptive measurement should satisfy? If we can find an adaptive measurement achieving the Helstrom limit, what is the relationship between this adaptive measure-

ment and the optimum joint measurement in the subspace spanned by M -input states, $\{|\psi_i^N\rangle\}$, $i = 1, \dots, M$? Answers for these questions are summarized in the following lemma, and it can be proved by directly applying Lemma 3.1.

Lemma 3.2. *When $\mathcal{P}_{\mathcal{U}_M}$ is a projector to the M -dimensional subspace \mathcal{U}_M spanned by $\{|\psi_i^N\rangle\}$, the projection of Π_i in (3.9) onto \mathcal{U}_M , i.e., $\Pi'_i = \mathcal{P}_{\mathcal{U}_M} \Pi_i \mathcal{P}_{\mathcal{U}_M}$, can achieve the Helstrom limit if and only if $\{\Pi'_i\}$ is an orthonormal projective measurement, i.e.,*

$$\Pi'_i \Pi'_j = \delta_{ij} \Pi'_i, \quad (3.10)$$

satisfying (3.3)–(3.4).

Proof. Lemma 3.1 shows that there exists a *unique* optimum measurement in \mathcal{U}_M , which is a set of orthonormal projectors. Therefore, for the adaptive measurement $\{\Pi_i\}$, operating on $\mathcal{H}^{\otimes N} \supset \mathcal{U}_M$, to perform as well as the optimum measurement and to achieve the Helstrom limit, the projection of $\{\Pi_i\}$ onto \mathcal{U}_M should be exactly the same as the optimum measurement in the subspace \mathcal{U}_M . \square

This lemma implies an important property of the optimum adaptive measurement $\{\Pi'_i\}$ summarized in the following lemma. Note that with adaptive measurement vectors $\{|\Omega(y_1^N)\rangle\}$, $y_1^N \in \{1, \dots, M\}^N$, the probability of observing the output sequence y_1^N , given the true hypothesis H_j is $p(Y_1^N = y_1^N | H_j) := \left| \langle \psi_j^N | \Omega(y_1^N) \rangle \right|^2$.

Lemma 3.3. *For the adaptive measurement vectors $\{|\Omega(y_1^N)\rangle\}$, $y_1^N \in \{1, \dots, M\}^N$, which form a POVM $\{\Pi_i\}$ after grouping by D_i , $i = 1, \dots, M$, to achieve the Helstrom limit, they should satisfy*

$$\frac{p(Y_1^N = y_1^N | H_j)}{p(Y_1^N = y_1^N | H_k)} = \frac{\sum_{y_1^N \in D_i} p(Y_1^N = y_1^N | H_j)}{\sum_{y_1^N \in D_i} p(Y_1^N = y_1^N | H_k)} \quad (3.11)$$

for all y_1^N 's belonging to the i -th decision set D_i , $i = 1, \dots, M$, and for every $j, k = 1, \dots, M$. When we denote the event $\hat{H}_i = \{y_1^N : y_1^N \in D_i\}$, by Bayes' rule, (3.11) can

also be written as

$$\frac{p(H_j|\widehat{H}_i)}{p(H_k|\widehat{H}_i)} = \frac{p(H_j|Y_1^N = y_1^N)}{p(H_k|Y_1^N = y_1^N)}, \quad \forall y_1^N \in D_i. \quad (3.12)$$

Moreover, (3.12) implies that for any $j \in \{1, \dots, M\}$, $p(H_j|y_1^N)$ is the same for $\forall y_1^N \in D_i$.

Remark 3.1. An important property of the optimum adaptive measurement, which achieves the Helstrom limit, is summarized in (3.12). One property of the optimum adaptive measurement is that the ratio between the posterior probabilities of any two hypotheses H_j and H_k is the same for every y_1^N that belongs to the same decision set D_i . In other words, it means that the optimal adaptive measurement should guarantee the same quality of decision for every output $y_1^N \in D_i$, i.e., for every y_1^N that belongs to the same decision set D_i , the probability to make a right guess should be the same.

Proof. Let us denote a basis of the space \mathcal{U}_M as $\{|v_i\rangle\}$, $i = 1, \dots, M$. We can represent the basis with linearly independent input states $\{|\psi_i^N\rangle\} \in \mathcal{U}_M$, $i = 1, \dots, M$ by Gram-Schmidt process. First, let $|v_1\rangle := |\psi_1^N\rangle$, and then let

$$|v_i\rangle := \frac{1}{\sqrt{1 - \sum_{j=1}^{i-1} |\langle \psi_i^N | v_j \rangle|^2}} \left(|\psi_i^N\rangle - \sum_{j=1}^{i-1} \langle v_j | \psi_i^N \rangle |v_j\rangle \right) \quad (3.13)$$

for $i = 2, \dots, M$. Then the projector $\mathcal{P}_{\mathcal{U}_M}$ onto \mathcal{U}_M can be written with the basis $\{|v_i\rangle\}$ as $\mathcal{P}_{\mathcal{U}_M} = \sum_{i=1}^M |v_i\rangle\langle v_i|$.

From Lemma 3.2, we know that it is necessary that $\Pi'_i = \mathcal{P}_{\mathcal{U}_M} \Pi_i \mathcal{P}_{\mathcal{U}_M}$ should be a rank-1 operator, i.e., $\Pi'_i = |\omega\rangle\langle\omega|$ for some $|\omega\rangle \in \mathcal{U}_M$. The vector $|\omega\rangle \in \mathcal{U}_M$ can be written in terms of the basis $\{|v_i\rangle\}$ of \mathcal{U}_M , as $|\omega\rangle = \sum_{i=1}^M c_i |v_i\rangle$ for some $c_i \in \mathbb{C}$. From

the definition of Π_i in (3.9) and $\mathcal{P}_{\mathcal{U}_M}$,

$$\begin{aligned}\Pi'_i &= \left(\sum_{j=1}^M |v_j\rangle\langle v_j| \right) \cdot \left(\sum_{y_1^N \in D_i} |\Omega(y_1^N)\rangle\langle\Omega(y_1^N)| \right) \cdot \left(\sum_{k=1}^M |v_k\rangle\langle v_k| \right) \\ &= \sum_{j,k=1}^M \left(\sum_{y_1^N \in D_i} \langle v_j | \Omega(y_1^N) \rangle \langle \Omega(y_1^N) | v_k \rangle \right) |v_j\rangle\langle v_k|.\end{aligned}\quad (3.14)$$

This Π'_i should be a rank-1 operator in \mathcal{U}_M , i.e. $\Pi'_i = |\omega\rangle\langle\omega|$ for $|\omega\rangle = \sum_{i=1}^M c_i |v_i\rangle$, which can also be written as

$$\Pi'_i = \left(\sum_j^M c_j |v_j\rangle \right) \left(\sum_k^M c_k^* \langle v_k| \right) = \sum_{j,k=1}^M c_j c_k^* |v_j\rangle\langle v_k|. \quad (3.15)$$

For (3.14) to satisfy the form of (3.15), the coefficients in (3.14) should also satisfy $|c_j|^2 \cdot |c_k|^2 = |c_j c_k^*|^2$, which is implied from (3.15), and thus the following equations should be satisfied

$$\begin{aligned}& \left(\sum_{y_1^N \in D_i} |\langle v_j | \Omega(y_1^N) \rangle|^2 \right) \left(\sum_{y_1^N \in D_i} |\langle v_k | \Omega(y_1^N) \rangle|^2 \right) \\ &= \left| \sum_{y_1^N \in D_i} \langle v_j | \Omega(y_1^N) \rangle \langle \Omega(y_1^N) | v_k \rangle \right|^2\end{aligned}\quad (3.16)$$

with $\{|\Omega(y_1^N)\rangle\}$ for $i, j, k \in \{1, \dots, M\}$.

From the Cauchy-Schwarz inequality, we know that the equality in (3.16) can be satisfied if and only if

$$\langle v_j | \Omega(y_1^N) \rangle = \gamma_{j,k} \langle v_k | \Omega(y_1^N) \rangle \quad (3.17)$$

for all $y_1^N \in D_i$ with some constant $\gamma_{j,k} \in \mathbb{C}$. Then, from (3.17) and the basis $\{|v_i\rangle\}$ defined in (3.13), we can show that the adaptive measurement achieving the Helstrom limit should satisfy (3.11). This concludes the proof. \square

In the next section, we will consider *binary* hypothesis testing between $\{|\psi_0^N\rangle, |\psi_1^N\rangle\}$ and present the optimal adaptive measurement satisfying Lemma 3.2, and the corresponding properties in Lemma 3.3.

■ 3.4 Optimal Adaptive Measurements for Binary Hypothesis Testing

In this section, we will use Lemma 3.2 to derive the properties that should be satisfied with the optimal adaptive measurement achieving the Helstrom limit for the BHT between $|\psi_0^N\rangle$ and $|\psi_1^N\rangle$ with prior probabilities $\{p_0, p_1\}$. We will show that a greedy algorithm, combined with posterior updating, gives an adaptive measurement satisfying these necessary and sufficient conditions. Moreover, it will be shown that the Dolinar receiver [10], which has been known to achieve the Helstrom limit for the BHT between two *coherent* states, indeed physically implements the optimum adaptive measurement by a simple receiver structure including a photon counter with feedback control signal that is added to the received quantum state.

First, let us derive the optimum measurement operator $\{\Pi'_0, \Pi'_1\}$ acting in the subspace \mathcal{U}_2 spanned by $\{|\psi_0^N\rangle, |\psi_1^N\rangle\}$. We can find a basis vector $\{|x\rangle, |y\rangle\}$ of \mathcal{U}_2 such that

$$|\psi_0^N\rangle = \cos\theta|x\rangle + \sin\theta|y\rangle; \quad |\psi_1^N\rangle = \cos\theta|x\rangle - \sin\theta|y\rangle. \quad (3.18)$$

From Lemma 3.1, the optimum measurement operators are orthonormal projectors so that when we denote $\Pi'_0 = |\Omega_0\rangle\langle\Omega_0|$ and $\Pi'_1 = |\Omega_1\rangle\langle\Omega_1|$, the measurement vectors $|\Omega_0\rangle$ and $|\Omega_1\rangle$ are orthonormal to each other and can be written with a parameter ϕ as

$$|\Omega_0\rangle = \cos\phi|x\rangle + \sin\phi|y\rangle; \quad |\Omega_1\rangle = \sin\phi|x\rangle - \cos\phi|y\rangle. \quad (3.19)$$

To minimize the probability of error

$$p_e = 1 - \sum_{i=0}^1 p_i |\langle\psi_i^N|\Omega_i\rangle|^2 = p_0 \sin^2(\phi - \theta) + p_1 \cos^2(\phi + \theta), \quad (3.20)$$

the optimum measurement vectors, or the parameter ϕ^* , should satisfy

$$\begin{aligned} \frac{p_1^2}{p_0^2} &= \frac{\cos^2(\phi^* - \theta) \sin^2(\phi^* - \theta)}{\cos^2(\phi^* + \theta) \sin^2(\phi^* + \theta)} \\ &= \frac{|\langle \psi_0^N | \Omega_0 \rangle|^2 \cdot |\langle \psi_0^N | \Omega_1 \rangle|^2}{|\langle \psi_1^N | \Omega_0 \rangle|^2 \cdot |\langle \psi_1^N | \Omega_1 \rangle|^2} = \frac{p(0|0) \cdot p(1|0)}{p(0|1) \cdot p(1|1)}. \end{aligned} \quad (3.21)$$

The resulting minimum probability of error is

$$p_e^* = \frac{1 - \sqrt{1 - 4p_0p_1|\langle \psi_0 | \psi_1 \rangle|^{2N}}}{2}. \quad (3.22)$$

For $0 < |\langle \psi_0 | \psi_1 \rangle| < 1$, as N increases, the Helstrom limit, p_e^* , can be approximated as

$$p_e^* \approx p_0p_1|\langle \psi_0 | \psi_1 \rangle|^{2N}. \quad (3.23)$$

Note that the optimum measurement vectors $\{|\Omega_0\rangle, |\Omega_1\rangle\}$ that achieve this Helstrom limit might be an entangling measurement in \mathcal{U}_2 .

■ 3.4.1 Necessary and Sufficient Conditions for Adaptive Measurement

Now we will derive the necessary and sufficient conditions for the adaptive measurement, by using Lemma 3.2, to achieve the Helstrom limit of the BHT between $|\psi_0^N\rangle$ and $|\psi_1^N\rangle$. We consider the adaptive measurement that observes each state one at a time and generates an output $y_n \in \{0, 1\}$ for $n = 1, \dots, N$ with the ability to update the $(n+1)$ -th measurement based on the previous observations, y_1^n . The grouping strategy for the resulting 2^N -output of the adaptive measurement, $y_1^N \in \{0, 1\}^N$, is as follows: the grouping is based on the last output y_N ; when $y_N = 0$, \widehat{H}_0 is claimed, and when $y_N = 1$, \widehat{H}_1 . Therefore, the measurement operator from this adaptive measurement is

$$\Pi_i = \sum_{y_1^{N-1} \in \{0,1\}^{N-1}} \left| \Omega(y_1^{N-1}i) \right\rangle \left\langle \Omega(y_1^{N-1}i) \right|, \quad i \in \{0, 1\}, \quad (3.24)$$

where

$$|\Omega(y_1^N)\rangle := |\omega_1(y_1)\rangle \otimes \cdots \otimes |\omega_N(y_1^N)\rangle \quad (3.25)$$

for $y_1^N \in \{0, 1\}^N$. From Lemma 3.2, for the adaptive measurement in (3.24) to achieve the Helstrom limit, the projection of Π_i onto \mathcal{U}_2 should be the orthonormal projective measurement satisfying the optimality condition in (3.21), i.e.,

$$\mathcal{P}_{\mathcal{U}_2} \Pi_i \mathcal{P}_{\mathcal{U}_2} = |\Omega_i\rangle\langle\Omega_i|, \quad i = 0, 1 \quad (3.26)$$

where $\{|\Omega_0\rangle, |\Omega_1\rangle\}$ are the optimal measurements in \mathcal{U}_2 , satisfying (3.21), and $\mathcal{P}_{\mathcal{U}_2}$ is the projector onto the subspace \mathcal{U}_2 spanned by $\{|\psi_0^N\rangle, |\psi_1^N\rangle\}$.

Since a basis $\{|v_1\rangle, |v_2\rangle\}$ of the 2-dimensional subspace \mathcal{U}_2 can be written in terms of $\{|\psi_0^N\rangle, |\psi_1^N\rangle\}$ as

$$|v_1\rangle = |\psi_0^N\rangle, \quad |v_2\rangle = \frac{1}{\sqrt{1 - |\langle\psi_0^N|\psi_1^N\rangle|^2}} (|\psi_1^N\rangle - (\langle\psi_0^N|\psi_1^N\rangle) |\psi_0^N\rangle), \quad (3.27)$$

we can write the projector $\mathcal{P}_{\mathcal{U}_2}$ with this basis, as $\mathcal{P}_{\mathcal{U}_2} = |v_1\rangle\langle v_1| + |v_2\rangle\langle v_2|$.

By plugging this $\mathcal{P}_{\mathcal{U}_2}$, which is written in terms of $\{|\psi_0^N\rangle, |\psi_1^N\rangle\}$, into (3.26), we can write down the necessary and sufficient conditions for the adaptive measurement vectors $\{|\Omega(y_1^N)\rangle\}$, $y_1^N \in \{0, 1\}^N$, explicitly in terms of the input states $\{|\psi_0^N\rangle, |\psi_1^N\rangle\}$. By using

$$\begin{aligned} p(Y_1^N = y_1^N | H_i) &= |\langle\psi_i^N|\Omega(y_1^N)\rangle|^2, \\ p(H_j | H_i) &= p(Y_N = j | H_i) = \sum_{y_1^N \in \{0,1\}^{N-1}} |\langle\psi_i^N|\Omega(y_1^{N-1}j)\rangle|^2, \end{aligned} \quad (3.28)$$

for $i, j = 0, 1$, we can now rephrase (3.26) in terms of the probability distributions, as summarized below.

Lemma 3.4. *The necessary and sufficient conditions for the adaptive measurement to*

achieve the Helstrom limit for the BHT between $|\psi_0^N\rangle$ and $|\psi_1^N\rangle$ are

$$\frac{p_0}{p_1} \cdot \frac{p(Y_1^N = y_1^{N-1}0, H_1)}{p(Y_1^N = y_1^{N-1}0, H_0)} = \left(c^N - \sqrt{(1 - c^{2N}) \frac{p(Y_N = 1, H_0)}{p(Y_N = 0, H_0)}} \right)^2, \quad (3.29)$$

$$\frac{p_0}{p_1} \cdot \frac{P(Y_1^N = y_1^{N-1}1, H_1)}{P(Y_1^N = y_1^{N-1}1, H_0)} = \left(c^N + \sqrt{(1 - c^{2N}) \frac{p(Y_N = 0, H_0)}{p(Y_N = 1, H_0)}} \right)^2, \quad (3.30)$$

$$\frac{p(Y_N = 0, H_0) \cdot p(Y_N = 1, H_0)}{p(Y_N = 0, H_1) \cdot p(Y_N = 1, H_1)} = 1, \quad (3.31)$$

for $\forall y_1^{N-1} \in \{0, 1\}^{N-1}$ where $c := |\langle \psi_0 | \psi_1 \rangle|$, and thus $c^N = |\langle \psi_0 | \psi_1 \rangle|^N = |\langle \psi_0^N | \psi_1^N \rangle|$.

Proof. Appendix 3.A. □

Remark 3.2. From (3.29) and (3.30), respectively, it can be inferred that the optimal adaptive measurement should satisfy

$$\begin{aligned} p(H_0 | Y_N = 0) &= p(H_0 | Y_1^N = y_1^{N-1}0), \\ p(H_1 | Y_N = 1) &= p(H_1 | Y_1^N = y_1^{N-1}1), \end{aligned} \quad (3.32)$$

for $\forall y_1^{N-1} \in \{0, 1\}^{N-1}$, meaning that for any trajectory of y_1^{N-1} , the quality of the final decision should be the same for every y_1^N that belongs to the same decision set. Moreover, (3.31) combined with (3.32) implies that the optimum adaptive measurement should satisfy

$$p(H_0 | Y_1^N = y_1^{N-1}0) = p(H_1 | Y_1^N = z_1^{N-1}1), \quad (3.33)$$

for $y_1^N, z_1^N \in \{0, 1\}^N$. In other words, regardless of the final decision $Y_N = 0, 1$ as well as the previous outputs y_1^{N-1}, z_1^{N-1} , the quality of the final decision (the probability of making the right guess, given a length- N output sequence) should be the same.

■ 3.4.2 Bayesian Updating Rules

Now the question is how to find the adaptive measurement vectors $\{|\Omega(y_1^N)\rangle\}$, $y_1^N \in \{0, 1\}^N$ satisfying (3.29)–(3.31), which are the necessary and sufficient conditions to

achieve the Helstrom limit. In [1], Acin *et al.* showed that a *greedy algorithm*, which makes the locally optimal choice at each stage without considering its consequent effect in the future stages, coupled with *the posterior updating*, results in the globally optimal solution, i.e., achieving the Helstrom limit.

Consider the r -th stage optimization for $r = 1, \dots, N$, where the previous observations from measurements are denoted as $y_1^{r-1} \in \{0, 1\}^{r-1}$. The posterior distributions for the two hypotheses H_0 and H_1 is written as $\{p(H_0|Y_1^{r-1} = y_1^{r-1}), p(H_1|Y_1^{r-1} = y_1^{r-1})\}$. When $r = 1$, we just consider the prior probabilities, i.e., $\{p_0, p_1\}$. Assume that we choose the adaptive measurement of the r -th stage's, $\{|\omega(y_1^{r-1}0)\rangle, |\omega(y_1^{r-1}1)\rangle\}$, that minimizes the current stage's probability of error from the view of hard-decision at the current stage, which is

$$\begin{aligned} p_e^{(r)}(y_1^{r-1}) = & p(H_0|Y_1^{r-1} = y_1^{r-1}) \cdot |\langle \psi_0 | \omega(y_1^{r-1}1) \rangle|^2 \\ & + p(H_1|Y_1^{r-1} = y_1^{r-1}) \cdot |\langle \psi_1 | \omega(y_1^{r-1}0) \rangle|^2. \end{aligned} \quad (3.34)$$

To minimize (3.34), $\{|\omega(y_1^r 0)\rangle, |\omega(y_1^r 1)\rangle\}$ should meet

$$\begin{aligned} \frac{p(H_1|Y_1^{r-1} = y_1^{r-1})^2}{p(H_0|Y_1^{r-1} = y_1^{r-1})^2} &= \frac{|\langle \psi_0 | \omega(y_1^{r-1}0) \rangle|^2 \cdot |\langle \psi_0 | \omega(y_1^{r-1}1) \rangle|^2}{|\langle \psi_1 | \omega(y_1^{r-1}0) \rangle|^2 \cdot |\langle \psi_1 | \omega(y_1^{r-1}1) \rangle|^2} \\ &= \frac{p(Y_1^r = y_1^{r-1}0|H_0) \cdot p(Y_1^r = y_1^{r-1}1|H_0)}{p(Y_1^r = y_1^{r-1}0|H_1) \cdot p(Y_1^r = y_1^{r-1}1|H_1)}. \end{aligned} \quad (3.35)$$

In [1], it was shown that this greedy strategy achieves the Helstrom limit in (3.22) after any N -stages.

Now we show that this strategy indeed satisfies the necessary and sufficient conditions for the optimum adaptive measurement in Lemma 3.4.

Lemma 3.5. *When the r -th adaptive measurement vectors $\{|\omega(y_1^{r-1}0)\rangle, |\omega(y_1^{r-1}1)\rangle\}$, $y_1^{r-1} \in \{0, 1\}^{r-1}$, are chosen to minimize the current stage's probability of error (3.34) for every $r = 1, \dots, N$, those vectors satisfy the necessary and sufficient conditions, (3.29)–(3.31), to achieve the Helstrom limit.*

Proof. Appendix 3.B. □

■ 3.4.3 Dolinar Receiver

In this section, we will consider the BHT between two *coherent states* $\{|\gamma_0\rangle, |\gamma_1\rangle\}$. A *coherent state* $|\gamma\rangle$ is the quantum description of a single spatio-temporal-polarization mode of a classical optical-frequency electromagnetic (ideal laser-light) field, where $\gamma \in \mathbb{C}$ is the complex amplitude, and $|\gamma|^2$ is the mean photon number of the mode. The coherent state can be represented as a linear combination of number states $|n\rangle$ as

$$|\gamma\rangle = e^{-\frac{|\gamma|^2}{2}} \sum_{n=0}^{\infty} \frac{\gamma^n}{\sqrt{n!}} |n\rangle. \quad (3.36)$$

Therefore, when we measure the quantum state by a set of (photon) number state, or Fock state, $\{|n\rangle\langle n|\}$, $n = 0, \dots, \infty$, which can be physically implemented by a photon counter, the output of the measurement follows a Poisson distribution with rate $|\gamma|^2$, i.e.,

$$\Pr(n\text{-photon arrival}) = \frac{|\gamma|^{2n} e^{-|\gamma|^2}}{n!}. \quad (3.37)$$

The magnitude squared of the inner product between two coherent states $\{|\gamma_0\rangle, |\gamma_1\rangle\}$ is $|\langle\gamma_0|\gamma_1\rangle|^2 = e^{-|\gamma_0-\gamma_1|^2}$, and the Helstrom limit (the minimum average detection error probability) for these coherent states is

$$p_e^* = \frac{1 - \sqrt{1 - 4p_0p_1|\langle\gamma_0|\gamma_1\rangle|^2}}{2} = \frac{1 - \sqrt{1 - 4p_0p_1e^{-|\gamma_0-\gamma_1|^2}}}{2}, \quad (3.38)$$

where the prior probabilities are $\{p_0, p_1\}$, respectively.

Due to a peculiar property of the coherent state, a coherent state of duration T can be represented as a tensor product of N -copy of a weaker coherent state of duration T/N as

$$|\gamma\rangle = \underbrace{\left| \frac{\gamma}{\sqrt{N}} \right\rangle \otimes \dots \otimes \left| \frac{\gamma}{\sqrt{N}} \right\rangle}_N. \quad (3.39)$$

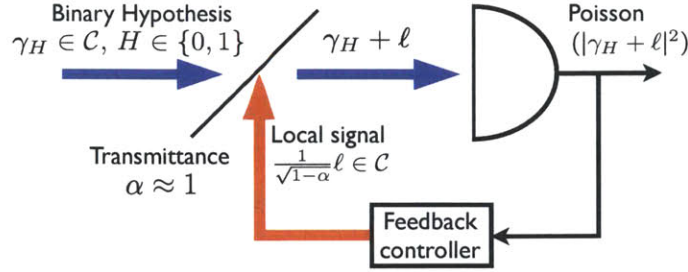


Figure 3.1. Coherent receiver using local feedback signal

Remember that the optimum adaptive measurement, satisfying (3.35), achieves the Helstrom bound for the BHT between *any* two length- N (identical-copy) tensor product states $\{|\psi_0^N\rangle, |\psi_1^N\rangle\}$. Therefore, we can apply the analysis in the previous section to the BHT between two coherent states. This was first observed in [2].

Even before this observation, which connects the BHT between N -copy of a quantum state with that of coherent states, Dolinar in [10] originally suggested a receiver (depicted in Fig. 3.1) that achieves the Helstrom bound in (3.38), by using the adaptive measurement techniques. Dolinar receiver basically combines a received coherent state with a controlled feedback coherent state, and then feed the merged signal to a photon counter. The feedback control signal can be adaptively updated based on the previous output observations. Dolinar proved that this receiver achieves the Helstrom bound with the feedback control signal optimized from the view of dynamic programming.

In the rest of this section, we will show that the Dolinar receiver works exactly as does the optimal adaptive measurement as described in (3.35), for two input states $\left\{ \left| \frac{\gamma_0}{\sqrt{N}} \right\rangle, \left| \frac{\gamma_1}{\sqrt{N}} \right\rangle \right\}$, in the limit of $N \rightarrow \infty$. When we denote $\delta := 1/N$, the two coherent states can be approximated from (3.36) as

$$\begin{aligned} \left| \frac{\gamma_0}{\sqrt{N}} \right\rangle &= \sqrt{1 - |\gamma_0|^2 \delta} |0\rangle + \gamma_0 \sqrt{\delta} |1\rangle + O(\delta) \\ \left| \frac{\gamma_1}{\sqrt{N}} \right\rangle &= \sqrt{1 - |\gamma_1|^2 \delta} |0\rangle + \gamma_1 \sqrt{\delta} |1\rangle + O(\delta) \end{aligned} \quad (3.40)$$

as $N \rightarrow \infty$. Thus the coherent state of weak energy can be approximated as a state

staying in the 2-dimensional space spanned by the two number states $\{|0\rangle, |1\rangle\}$. It means that when we measure such a weak coherent state, with a high probability, we observe either 0 or 1 photon.

Two orthonormal measurement vectors for these weak coherent states can also stay in the same 2-dimensional space and can be written as

$$\begin{aligned} |\omega_0\rangle &= \sqrt{1-\beta^2}|0\rangle + \beta e^{i\theta}|1\rangle, \\ |\omega_1\rangle &= \beta|0\rangle - e^{i\theta}\sqrt{1-\beta^2}|1\rangle \end{aligned} \quad (3.41)$$

with the basis $\{|0\rangle, |1\rangle\}$ for two parameters $\beta, \theta \in \mathbb{R}$.

To satisfy the optimality condition of the r -th adaptive measurement in (3.35), the parameters (β, θ) for the r -th adaptive measurement $\{|\omega(y_1^{r-1}0)\rangle, |\omega(y_1^{r-1}1)\rangle\}$ should be

$$e^{i2\theta} = \frac{\pi_0\gamma_0\sqrt{1-|\gamma_1|^2\delta} - \pi_1\gamma_1\sqrt{1-|\gamma_1|^2\delta}}{\pi_0\gamma_0^*\sqrt{1-|\gamma_1|^2\delta} - \pi_1\gamma_1^*\sqrt{1-|\gamma_1|^2\delta}}, \quad (3.42)$$

$$(\beta)^2 = \begin{cases} \frac{1}{2} \left(1 - \frac{|x|}{\sqrt{x^2+4}}\right) & \text{if } \pi_0 \geq \pi_1, \\ \frac{1}{2} \left(1 + \frac{|x|}{\sqrt{x^2+4}}\right) & \text{o.w.} \end{cases} \quad (3.43)$$

where

$$x = \frac{1}{\sqrt{\delta}} \left(\frac{\pi_0 - \pi_1}{|\pi_0\gamma_0 - \pi_1\gamma_1|} + O(\delta) \right), \quad (3.44)$$

and π_0 and π_1 are posterior probabilities at the r -th step, given previous observations $y_1^{r-1} \in \{0, 1\}^{r-1}$, i.e., $\pi_i = p(H_i | Y_1^{r-1} = y_1^{r-1})$ for $i = 0, 1$.

Then, how can we translate this (mathematically described) optimum adaptive measurement into a physical realization? The idea of the Dolinar receiver is to *rotate* two input states $\left\{ \left| \frac{\gamma_0}{\sqrt{N}} \right\rangle, \left| \frac{\gamma_1}{\sqrt{N}} \right\rangle \right\}$ in the subspace spanned by $\{|0\rangle, |1\rangle\}$ to make the optimal measurement of the rotated input states be exactly the number states $\{|0\rangle, |1\rangle\}$, which can be realized by a photon counter. The rotation of the input states $\left\{ \left| \frac{\gamma_0}{\sqrt{N}} \right\rangle, \left| \frac{\gamma_1}{\sqrt{N}} \right\rangle \right\}$ can be realized by a displacement operator, i.e., adding coherent state $\left| \frac{\ell}{\sqrt{N}} \right\rangle$ to the

input states $\left| \frac{\gamma_i}{\sqrt{N}} \right\rangle$, $i = 0, 1$, making it $\left| \frac{\gamma_i + \ell}{\sqrt{N}} \right\rangle$. Note that this process preserves the inner product between the two input states before and after the transformation, i.e.,

$$\left| \left\langle \frac{\gamma_0}{\sqrt{N}} \middle| \frac{\gamma_1}{\sqrt{N}} \right\rangle \right|^2 = \exp(-|\gamma_0 - \gamma_1|^2) = \left| \left\langle \frac{\gamma_0 + \ell}{\sqrt{N}} \middle| \frac{\gamma_1 + \ell}{\sqrt{N}} \right\rangle \right|^2. \quad (3.45)$$

Then the question is how to design ℓ to make the optimal measurement for the merged input states $\left\{ \left| \frac{\gamma_0 + \ell}{\sqrt{N}} \right\rangle, \left| \frac{\gamma_1 + \ell}{\sqrt{N}} \right\rangle \right\}$ become $\{|0\rangle, |1\rangle\}$. From (3.41) and (3.43), we can see that when $\beta = 0$, the optimum measurement becomes $|\omega_0\rangle = |0\rangle$ and $|\omega_1\rangle = |1\rangle$. To make $\beta = 0$, from (3.43), we need to make $|x| \rightarrow \infty$. When we replace γ_0 and γ_1 by $\gamma_0 + \ell$ and $\gamma_1 + \ell$, it can be shown that when

$$\ell = -\frac{\pi_0 \gamma_0 - \pi_1 \gamma_1}{\pi_0 - \pi_1}, \quad (3.46)$$

the denominator of x in (3.44) becomes 0, which makes $x \rightarrow \infty$. In other words, the optimum adaptive measurement at the r -th step, can be implemented by the Dolinar receiver that updates its control signal ℓ to be (3.46) where $\{\pi_0, \pi_1\} = \{p(H_0|Y_1^{r-1} = y_1^{r-1}), p(H_1|Y_1^{r-1} = y_1^{r-1})\}$.

Our solution for the optimum feedback control signal of the Dolinar receiver, which is derived from the view of optimum adaptive measurement, coincides with that derived in the original paper [10] by using much more complicated dynamic programmings.

■ 3.5 Conclusion

We considered the M -ary hypothesis testing problem among nonorthogonal quantum states, when N -copy of the unknown quantum state is available. The most general form of the measurement, including *entangling measurement*, can achieve the theoretical lower bound of the detection error probability (Helstrom limit). But since it is hard to physically implement the entangling measurement, our focus was on the performance of the adaptive (product) measurement, with which we measure one system at a time and then update the next measurement based on the previous observations.

We derived the necessary and sufficient conditions for the adaptive measurement to perform as well as the optimal entangling measurement that achieves the Helstrom limit. We then considered the binary hypothesis testing (BHT) problem, and showed that a greedy algorithm with posterior updating can meet the necessary and sufficient conditions for the optimum adaptive measurement. Moreover, we showed that the Dolinar receiver, which has been known to perform optimally for the BHT between two coherent states, is an exact physical translation of the optimal adaptive measurement.

■ 3.A Proof of Lemma 3.4

The adaptive measurement for the binary hypothesis testing can be written as

$$\begin{aligned}\Pi_0 &= \sum_{y_1^{N-1} \in \{0,1\}^{N-1}} \left| \Omega(y_1^{N-1}0) \right\rangle \left\langle \Omega(y_1^{N-1}0) \right|, \\ \Pi_1 &= \sum_{y_1^{N-1} \in \{0,1\}^{N-1}} \left| \Omega(y_1^{N-1}1) \right\rangle \left\langle \Omega(y_1^{N-1}1) \right|,\end{aligned}\tag{3.47}$$

where the final decision happens depending on the last output y_N .

When we denote the optimum measurement for $\{|\psi_0^N\rangle, |\psi_1^N\rangle\}$ on \mathcal{U}_2 as $\{|\Omega_0\rangle, |\Omega_1\rangle\} \in \mathcal{U}_2$, the uniqueness of the optimum measurement on the subspace \mathcal{U}_2 , stated in Lemma 3.1, implies that the measurement operators $\{\Pi_0, \Pi_1\}$ can achieve the Helstrom limit if and only if

$$|\Omega_0\rangle\langle\Omega_0| = P_{\mathcal{U}_2}\Pi_0P_{\mathcal{U}_2}, \quad |\Omega_1\rangle\langle\Omega_1| = P_{\mathcal{U}_2}\Pi_1P_{\mathcal{U}_2},\tag{3.48}$$

when $P_{\mathcal{U}_2}$ is the projective operator onto \mathcal{U}_2 . These conditions will be used to derive more specified necessary and sufficient conditions for the adaptive measurements to achieve the Helstrom limit.

When an orthonormal basis for the 2-dimensional space \mathcal{U}_2 is denoted as $\{|v_0\rangle, |v_1\rangle\}$,

the conditions in (3.48) can be written as

$$\begin{aligned}
|\Omega_0\rangle\langle\Omega_0| &= (|v_0\rangle\langle v_0| + |v_1\rangle\langle v_1|)\Pi_0(|v_0\rangle\langle v_0| + |v_1\rangle\langle v_1|) \\
&= \sum_{i=0}^{2^{N-1}-1} \left(|a_i|^2 |v_0\rangle\langle v_0| + a_i^* b_i |v_0\rangle\langle v_1| + b_i^* a_i |v_1\rangle\langle v_0| + |b_i|^2 |v_1\rangle\langle v_1| \right) \\
|\Omega_1\rangle\langle\Omega_1| &= (|v_0\rangle\langle v_0| + |v_1\rangle\langle v_1|)\Pi_1(|v_0\rangle\langle v_0| + |v_1\rangle\langle v_1|) \\
&= \sum_{i=0}^{2^{N-1}-1} \left(|c_i|^2 |v_0\rangle\langle v_0| + c_i^* d_i |v_0\rangle\langle v_1| + d_i^* c_i |v_1\rangle\langle v_0| + |d_i|^2 |v_1\rangle\langle v_1| \right)
\end{aligned} \tag{3.49}$$

where

$$\begin{aligned}
a_i &= \langle \Omega(y_1^{N-1}0) | v_0 \rangle; & b_i &= \langle \Omega(y_1^{N-1}0) | v_1 \rangle, \\
c_i &= \langle \Omega(y_1^{N-1}1) | v_0 \rangle; & d_i &= \langle \Omega(y_1^{N-1}1) | v_1 \rangle,
\end{aligned} \tag{3.50}$$

for the index $i = \sum_{j=1}^{N-1} y_j \cdot 2^{j-1}$ for each $y_1^{N-1} \in \{0, 1\}^{N-1}$.

The optimum measurement vectors $\{|\Omega_0\rangle, |\Omega_1\rangle\} \in \mathcal{U}_2$ can be written

$$\begin{aligned}
|\Omega_0\rangle &= \sqrt{1-c^2}|v_0\rangle + ce^{i\theta}|v_1\rangle \\
|\Omega_1\rangle &= c|v_0\rangle - e^{i\theta}\sqrt{1-c^2}|v_1\rangle
\end{aligned} \tag{3.51}$$

for $c \geq 0$ and $\theta \in \mathbb{R}$, and thus

$$\begin{aligned}
|\Omega_0\rangle\langle\Omega_0| &= (1-c^2)|v_0\rangle\langle v_0| + c\sqrt{1-c^2}e^{-i\theta}|v_0\rangle\langle v_1| + c\sqrt{1-c^2}e^{i\theta}|v_1\rangle\langle v_0| + c^2|v_1\rangle\langle v_1| \\
|\Omega_1\rangle\langle\Omega_1| &= c^2|v_0\rangle\langle v_0| - c\sqrt{1-c^2}e^{-i\theta}|v_0\rangle\langle v_1| - c\sqrt{1-c^2}e^{i\theta}|v_1\rangle\langle v_0| + (1-c^2)|v_1\rangle\langle v_1|.
\end{aligned} \tag{3.52}$$

These optimum measurement vectors should satisfy (3.21), which is written again below,

$$\frac{p_1^2}{p_0^2} = \frac{|\langle\psi_0^N|\Omega_0\rangle|^2 \cdot |\langle\psi_0^N|\Omega_1\rangle|^2}{|\langle\psi_1^N|\Omega_0\rangle|^2 \cdot |\langle\psi_1^N|\Omega_1\rangle|^2}, \tag{3.53}$$

and this optimality equation specifies the parameters c and θ in (3.51).

The set of adaptive measurement vectors $\{|\Omega(y_1^N)\rangle\}$, $y_1^N \in \{0, 1\}^N$ can satisfy the necessary and sufficient conditions in (3.49), if and only if the following three conditions can be met.

- $P_{\mathcal{U}_2}\Pi_0P_{\mathcal{U}_2}$ and $P_{\mathcal{U}_2}\Pi_1P_{\mathcal{U}_2}$ are rank-1 operators, i.e., we can find $|\Omega'_0\rangle, |\Omega'_1\rangle$ such that $P_{\mathcal{U}_2}\Pi_0P_{\mathcal{U}_2} = |\Omega'_0\rangle\langle\Omega'_0|$ and $P_{\mathcal{U}_2}\Pi_1P_{\mathcal{U}_2} = |\Omega'_1\rangle\langle\Omega'_1|$ where $\{|\Omega'_0\rangle, |\Omega'_1\rangle\} \in \mathcal{U}_2$.
- The measurement vectors $|\Omega'_0\rangle, |\Omega'_1\rangle$ are orthonormal to each other.
- Moreover, $|\Omega'_0\rangle = |\Omega_0\rangle$, $|\Omega'_1\rangle = |\Omega_1\rangle$, i.e., it satisfies (3.53).

These conditions specify the coefficients a_i, b_i, c_i , and d_i 's in (3.50). Let us consider the first condition that $P_{\mathcal{U}_2}\Pi_0P_{\mathcal{U}_2}$ and $P_{\mathcal{U}_2}\Pi_1P_{\mathcal{U}_2}$ should be rank-1 operators. The operators are rank-1 if and only if the coefficients a_i, b_i, c_i , and d_i satisfy

$$\begin{aligned} \left(\sum_{i=0}^{2^N-1} |a_i|^2\right) \left(\sum_{i=0}^{2^N-1} |b_i|^2\right) &= \left(\sum_{i=0}^{2^N-1} a_i^* b_i\right) \left(\sum_{i=0}^{2^N-1} a_i b_i^*\right) = \left|\sum_{i=0}^{2^N-1} a_i b_i^*\right|^2, \\ \left(\sum_{i=0}^{2^N-1} |c_i|^2\right) \left(\sum_{i=0}^{2^N-1} |d_i|^2\right) &= \left(\sum_{i=0}^{2^N-1} c_i^* d_i\right) \left(\sum_{i=0}^{2^N-1} c_i d_i^*\right) = \left|\sum_{i=0}^{2^N-1} c_i d_i^*\right|^2. \end{aligned} \quad (3.54)$$

From Cauchy-Schwarz inequality, for $x_1, x_2, \dots, x_n \in \mathbb{C}$ and $y_1, y_2, \dots, y_n \in \mathbb{C}$,

$$\left|\sum_{i=1}^n x_i y_i^*\right|^2 \leq \left(\sum_{j=1}^n |x_j|^2\right) \left(\sum_{k=1}^n |y_k|^2\right), \quad (3.55)$$

while equality is achieved if and only if x and y are linearly dependent, i.e., $x_i = \gamma \cdot y_i$, $\forall i$, for some $\gamma \in \mathbb{C}$. Therefore, (3.54) can be satisfied if and only if

$$b_i = \gamma \cdot a_i \quad \text{and} \quad d_i = \beta \cdot c_i, \quad \forall i \in \{0, \dots, 2^N-1\} \quad (3.56)$$

for some $\gamma, \beta \in \mathbb{C}$.

Under (3.56), the projected operators $P_{\mathcal{U}_2}\Pi_0P_{\mathcal{U}_2}$ and $P_{\mathcal{U}_2}\Pi_1P_{\mathcal{U}_2}$ can be written as

rank-1 operators as follows.

$$\begin{aligned}
P_{U_2} \Pi_0 P_{U_2} &= \sum_{i=0}^{2^{N-1}-1} \left(|a_i|^2 |v_0\rangle\langle v_0| + \gamma |a_i|^2 |v_0\rangle\langle v_1| + \gamma^* |a_i|^2 |v_1\rangle\langle v_0| + |\gamma|^2 |a_i|^2 |v_1\rangle\langle v_1| \right) \\
&= \left(\sqrt{\sum_i |a_i|^2} |v_0\rangle + \gamma^* \sqrt{\sum_i |a_i|^2} |v_1\rangle \right) \left(\sqrt{\sum_i |a_i|^2} \langle v_0| + \gamma \sqrt{\sum_i |a_i|^2} \langle v_1| \right), \\
P_{U_2} \Pi_1 P_{U_2} &= \sum_{i=0}^{2^{N-1}-1} \left(|c_i|^2 |v_0\rangle\langle v_0| + \beta |c_i|^2 |v_0\rangle\langle v_1| + \beta^* |c_i|^2 |v_1\rangle\langle v_0| + |\beta|^2 |c_i|^2 |v_1\rangle\langle v_1| \right) \\
&= \left(\sqrt{\sum_i |c_i|^2} |v_0\rangle + \beta^* \sqrt{\sum_i |c_i|^2} |v_1\rangle \right) \left(\sqrt{\sum_i |c_i|^2} \langle v_0| + \beta \sqrt{\sum_i |c_i|^2} \langle v_1| \right),
\end{aligned} \tag{3.57}$$

and thus

$$\begin{aligned}
|\Omega'_0\rangle &= \sqrt{\sum_i |a_i|^2} |v_0\rangle + \gamma^* \sqrt{\sum_i |a_i|^2} |v_1\rangle, \\
|\Omega'_1\rangle &= \sqrt{\sum_i |c_i|^2} |v_0\rangle + \beta^* \sqrt{\sum_i |c_i|^2} |v_1\rangle.
\end{aligned} \tag{3.58}$$

Then we put additional conditions on $a_i, b_i, c_i,$ and d_i to make $|\Omega'_0\rangle, |\Omega'_1\rangle$ to be orthonormal to each other. It means that those coefficients should satisfy

$$\begin{aligned}
(1 + |\gamma|^2) \left(\sum_i |a_i|^2 \right) &= (1 + |\beta|^2) \left(\sum_i |c_i|^2 \right) = 1, \\
(1 + \gamma\beta^*) \sqrt{\sum_i |a_i|^2} \sqrt{\sum_i |c_i|^2} &= 0.
\end{aligned} \tag{3.59}$$

Assuming $\sum_i |a_i|^2 \neq 0$ and $\sum_i |c_i|^2 \neq 0$, it implies that

$$\begin{aligned}
|\gamma|^2 &= 1 / \left(\sum_i |a_i|^2 \right) - 1, \quad |\beta|^2 = 1 / \left(\sum_i |c_i|^2 \right) - 1, \\
\gamma &= -1/\beta^*.
\end{aligned} \tag{3.60}$$

By using the fact that $\sum_i |a_i|^2 + \sum_j |c_j|^2 = 1$, these conditions can be written as

$$\gamma = -\sqrt{\frac{1 - (\sum_i |a_i|^2)}{(\sum_i |a_i|^2)}}, \quad \beta = \sqrt{\frac{(\sum_i |a_i|^2)}{1 - (\sum_i |a_i|^2)}}. \quad (3.61)$$

Now the last condition that a_i, b_i, c_i , and $d_i, \forall i$ should satisfy is for the resulting measurement vectors $\{|\Omega'_0\rangle, |\Omega'_1\rangle\} \in \mathcal{U}_2$, which can be written as

$$\begin{aligned} |\Omega'_0\rangle &= \left(\sqrt{\sum_i |a_i|^2} |v_0\rangle + \gamma^* \sqrt{\sum_i |a_i|^2} |v_1\rangle \right) = \left(\sqrt{\sum_i |a_i|^2} |v_0\rangle - \sqrt{1 - \sum_i |a_i|^2} |v_1\rangle \right), \\ |\Omega'_1\rangle &= \left(\sqrt{1 - \sum_i |a_i|^2} |v_0\rangle + \beta^* \sqrt{\sum_i |a_i|^2} |v_1\rangle \right) \\ &= \left(\sqrt{1 - \sum_i |a_i|^2} |v_0\rangle + \sqrt{\sum_i |a_i|^2} |v_1\rangle \right), \end{aligned} \quad (3.62)$$

the optimality condition (3.53) should hold.

In order to write that optimality condition in terms of the input states $\{|\Psi_0\rangle := |\psi_0^N\rangle, |\Psi_1\rangle := |\psi_1^N\rangle\}$, we first pick an orthonormal basis of \mathcal{U}_2 as follows.

$$|v_0\rangle = |\Psi_0\rangle, \quad |v_1\rangle = \frac{1}{\sqrt{1 - |\langle \Psi_0 | \Psi_1 \rangle|^2}} (|\Psi_1\rangle - \langle \Psi_0 | \Psi_1 \rangle |\Psi_0\rangle) \quad (3.63)$$

For such a basis

$$\sum_i |a_i|^2 = \sum_{y_1^{N-1}} \left| \langle \Omega(y_1^{N-1}0) | \Psi_0 \rangle \right|^2 = p(0|0), \quad \sum_i |c_i|^2 = 1 - \sum_i |a_i|^2 = p(1|0). \quad (3.64)$$

Moreover, the following inner product relationships hold.

$$\begin{aligned} \langle \Psi_0 | v_0 \rangle &= 1, \quad \langle \Psi_0 | v_1 \rangle = 0; \\ \langle \Psi_1 | v_0 \rangle &= \langle \Psi_1 | \Psi_0 \rangle, \quad \langle \Psi_1 | v_1 \rangle = \sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2}. \end{aligned} \quad (3.65)$$

Using these relationships and

$$|\Omega'_0\rangle = \sqrt{p(0|0)}|x\rangle - \sqrt{p(1|0)}|y\rangle, \quad |\Omega'_1\rangle = \sqrt{p(1|0)}|x\rangle + \sqrt{p(0|0)}|y\rangle, \quad (3.66)$$

it can be shown that

$$\begin{aligned} \langle \Psi_0 | \Omega'_0 \rangle &= \sqrt{p(0|0)} \langle \Psi_0 | v_0 \rangle - \sqrt{p(1|0)} \langle \Psi_0 | v_1 \rangle = \sqrt{p(0|0)}, \\ \langle \Psi_0 | \Omega'_1 \rangle &= \sqrt{p(1|0)} \langle \Psi_0 | v_0 \rangle + \sqrt{p(0|0)} \langle \Psi_0 | v_1 \rangle = \sqrt{p(1|0)}, \\ \langle \Psi_1 | \Omega'_0 \rangle &= \sqrt{p(0|0)} \langle \Psi_1 | v_0 \rangle - \sqrt{p(1|0)} \langle \Psi_1 | v_1 \rangle \\ &= \sqrt{p(0|0)} \langle \Psi_1 | \Psi_0 \rangle - \sqrt{p(1|0)} \sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2}, \\ \langle \Psi_1 | \Omega'_1 \rangle &= \sqrt{p(1|0)} \langle \Psi_1 | v_0 \rangle + \sqrt{p(0|0)} \langle \Psi_1 | v_1 \rangle \\ &= \sqrt{p(1|0)} \langle \Psi_1 | \Psi_0 \rangle + \sqrt{p(0|0)} \sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2}. \end{aligned} \quad (3.67)$$

Now we will show that $\langle \Psi_1 | \Omega'_0 \rangle$ is equal to $\sqrt{p(0|1)}$, which is

$$\sqrt{p(0|1)} = \sqrt{\sum_{y_1^{N-1}} |\langle \Omega(y_1^{N-1}0) | \Psi_1 \rangle|^2}. \quad (3.68)$$

From the definition of b_i and the basis $|v_1\rangle$,

$$b_i = \langle \Omega(y_1^{N-1}0) | v_1 \rangle = \frac{\langle \Omega(y_1^{N-1}0) | \Psi_1 \rangle - \langle \Psi_0 | \Psi_1 \rangle \langle \Omega(y_1^{N-1}0) | \Psi_0 \rangle}{\sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2}}, \quad (3.69)$$

which implies that

$$\langle \Omega(y_1^{N-1}0) | \Psi_1 \rangle = \langle \Psi_0 | \Psi_1 \rangle \langle \Omega(y_1^{N-1}0) | \Psi_0 \rangle + \sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2} \cdot b_i. \quad (3.70)$$

From the definition of a_i , and $b_i = \gamma \cdot a_i = -\sqrt{\frac{p(1|0)}{p(0|0)}} a_i$,

$$\langle \Omega(y_1^{N-1}0) | \Psi_1 \rangle = \left(\langle \Psi_0 | \Psi_1 \rangle - \sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2} \sqrt{\frac{p(1|0)}{p(0|0)}} \right) a_i. \quad (3.71)$$

When we sum up the magnitude squared of this term over all y_1^{N-1}

$$\begin{aligned}
\sum_{y_1^{N-1}} \left| \langle \Omega(y_1^{N-1}0) | \Psi_1 \rangle \right|^2 &= \left| \langle \Psi_0 | \Psi_1 \rangle - \sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2} \sqrt{\frac{p(1|0)}{p(0|0)}} \right|^2 \sum_i |a_i|^2 \\
&= \left| \langle \Psi_0 | \Psi_1 \rangle - \sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2} \sqrt{\frac{p(1|0)}{p(0|0)}} \right|^2 p(0|0) \\
&= \left| \sqrt{p(0|0)} \langle \Psi_0 | \Psi_1 \rangle - \sqrt{p(1|0)} \sqrt{1 - |\langle \Psi_1 | \Psi_0 \rangle|^2} \right|^2.
\end{aligned} \tag{3.72}$$

Therefore, by using this relationship and (3.67), it is shown that

$$|\langle \Psi_1 | \Omega'_0 \rangle|^2 = \sum_{y_1^{N-1}} \left| \langle \Omega(y_1^{N-1}0) | \Psi_1 \rangle \right|^2 = p(0|1). \tag{3.73}$$

Moreover, from $|\langle \Psi_1 | \Omega'_0 \rangle|^2 + |\langle \Psi_1 | \Omega'_1 \rangle|^2 = 1$,

$$|\langle \Psi_1 | \Omega'_1 \rangle|^2 = p(1|1). \tag{3.74}$$

By using (3.67), (3.73) and (3.74), the optimality condition in (3.53) can be written as

$$\frac{p_1^2}{p_0^2} = \frac{|\langle \Psi_0 | \Omega'_0 \rangle|^2 \cdot |\langle \Psi_0 | \Omega'_1 \rangle|^2}{|\langle \Psi_1 | \Omega'_0 \rangle|^2 \cdot |\langle \Psi_1 | \Omega'_1 \rangle|^2} = \frac{p(0|0) \cdot p(1|0)}{p(0|1) \cdot p(1|1)}. \tag{3.75}$$

In summary, the necessary and sufficient conditions for the adaptive measurements to achieve the Helstrom limit can be transformed to

- $b_i = \gamma a_i, \quad d_i = \beta c_i, \quad \forall i \{0, \dots, 2^{N-1} - 1\}$.
- $\gamma = -\sqrt{p(1|0)/p(0|0)}, \quad \beta = \sqrt{p(0|0)/p(1|0)}$ for basis $|v_0\rangle$ and $|v_1\rangle$ in (3.63)
- $p_1^2/p_0^2 = (p(0|0) \cdot p(1|0)) / (p(0|1) \cdot p(1|1))$

When we combine the first and second conditions, for basis $|v_0\rangle$ and $|v_1\rangle$ in (3.63),

$$\begin{aligned}\langle \Omega(y_1^{N-1}0) | \Psi_1 \rangle &= \left(\langle \Psi_0 | \Psi_1 \rangle - \sqrt{(1 - |\langle \Psi_0 | \Psi_1 \rangle|^2) \frac{p(1|0)}{p(0|0)}} \right) \langle \Omega(y_1^{N-1}0) | \Psi_0 \rangle \\ \langle \Omega(y_1^{N-1}1) | \Psi_1 \rangle &= \left(\langle \Psi_0 | \Psi_1 \rangle + \sqrt{(1 - |\langle \Psi_0 | \Psi_1 \rangle|^2) \frac{p(0|0)}{p(1|0)}} \right) \langle \Omega(y_1^{N-1}1) | \Psi_0 \rangle\end{aligned}\quad (3.76)$$

for all $y_1^{N-1} \in \{0, 1\}^{N-1}$.

Therefore, the necessary and sufficient conditions for the adaptive measurements can be summarized as in (3.29)-(3.31), i.e.,

$$\begin{aligned}\frac{p_0}{p_1} \cdot \frac{P(Y_1^N = y_1^{N-1}0, H_1)}{P(Y_1^N = y_1^{N-1}0, H_0)} &= \left(c^N - \sqrt{(1 - c^{2N}) \frac{P(Y_N = 1, H_0)}{P(Y_N = 0, H_0)}} \right)^2, \\ \frac{p_0}{p_1} \cdot \frac{P(Y_1^N = y_1^{N-1}1, H_1)}{P(Y_1^N = y_1^{N-1}1, H_0)} &= \left(c^N + \sqrt{(1 - c^{2N}) \frac{P(Y_N = 0, H_0)}{P(Y_N = 1, H_0)}} \right)^2, \\ \frac{P(Y_N = 0, H_0) \cdot P(Y_N = 1, H_0)}{P(Y_N = 0, H_1) \cdot P(Y_N = 1, H_1)} &= 1,\end{aligned}$$

for $\forall y_1^{N-1} \in \{0, 1\}^{N-1}$ where $c := |\langle \psi_0 | \psi_1 \rangle|$, and thus $c^N = |\langle \psi_0 | \psi_1 \rangle|^N = |\langle \psi_0^N | \psi_1^N \rangle|$.

■ 3.B Proof of Lemma 3.5

The adaptive measurement suggested in [1] follows a greedy algorithm paired with the posterior updating, and it satisfies

$$\frac{(p(H_1 | Y_1^{r-1} = y_1^{r-1}))^2}{(p(H_0 | Y_1^{r-1} = y_1^{r-1}))^2} = \frac{p(Y_r = 0 | H_0, Y_1^{r-1} = y_1^{r-1}) \cdot p(Y_r = 1 | H_0, Y_1^{r-1} = y_1^{r-1})}{p(Y_r = 0 | H_1, Y_1^{r-1} = y_1^{r-1}) \cdot p(Y_r = 1 | H_1, Y_1^{r-1} = y_1^{r-1})}\quad (3.77)$$

for every $r = 1, \dots, N$ and $\forall y_1^{r-1} \in \{0, 1\}^{r-1}$. From this updating rule, it can be shown that

$$p(Y_1^r = y_1^{r-1}a, H_b) = \frac{p(Y_1^{r-1} = y_1^{r-1}, H_b)}{2} \left(1 + (-1)^{a+b} \times \frac{p(Y_1^{r-1} = y_1^{r-1}, H_b) + (1 - 2c^2)p(Y_1^{r-1} = y_1^{r-1}, H_{b \oplus 1})}{R(y_1^{r-1})} \right) \quad (3.78)$$

where $a, b \in \{0, 1\}$ where $c := |\langle \psi_0 | \psi_1 \rangle|$. As stated in [1], (3.78) also implies that

$$p_0 p_1 c^{2r} [p(Y_1^r = y_1^r, H_0) + p(Y_1^r = y_1^r, H_1)]^2 - p(Y_1^r = y_1^r, H_0) \cdot p(Y_1^r = y_1^r, H_1) = 0 \quad (3.79)$$

for any $y_1^r \in \{0, 1\}^r$. We can show it by mathematical induction. First, we can check that (3.79) is true for $r = 1$, since

$$\begin{aligned} p(Y_1 = a, H_0) &= \frac{p_0}{2} \left(1 + (-1)^a \frac{p_0 + (1 - 2c^2)p_1}{\sqrt{1 - 4p_0p_1c^2}} \right), \\ p(Y_1 = a, H_1) &= \frac{p_1}{2} \left(1 - (-1)^a \frac{p_1 + (1 - 2c^2)p_0}{\sqrt{1 - 4p_0p_1c^2}} \right) \end{aligned} \quad (3.80)$$

for $a \in \{0, 1\}$, from (3.78). Assume that (3.79) is true for r , and then show that it is also true for $(r + 1)$. From (3.78),

$$\begin{aligned} p(Y_1^r a, H_0) &= \frac{p(Y_1^r = y_1^r, H_0)}{2} \cdot \frac{(-1)^a (1 - 2p_0p_1c^{2(r+1)}) + \sqrt{1 - 4p_0p_1c^{2(r+1)}}}{\sqrt{1 - 4p_0p_1c^{2(r+1)}}} \\ &\quad + \frac{p(Y_1^r = y_1^r, H_1)}{2} \cdot \frac{-(-1)^a 2p_0p_1c^{2(r+1)}}{\sqrt{1 - 4p_0p_1c^{2(r+1)}}}, \\ p(Y_1^r a, H_1) &= \frac{p(Y_1^r = y_1^r, H_0)}{2} \cdot \frac{(-1)^a 2p_0p_1c^{2(r+1)}}{\sqrt{1 - 4p_0p_1c^{2(r+1)}}} \\ &\quad + \frac{p(Y_1^r = y_1^r, H_1)}{2} \cdot \frac{-(-1)^a (1 - 2p_0p_1c^{2(r+1)}) + \sqrt{1 - 4p_0p_1c^{2(r+1)}}}{\sqrt{1 - 4p_0p_1c^{2(r+1)}}}. \end{aligned} \quad (3.81)$$

From this, we can directly show (3.79) is true for $(r + 1)$. Therefore, (3.79) is true for every integer r .

We will also show that (3.78) is equivalent to (3.77). Since it is straightforward to derive (3.78) from (3.77), here we will only show the reverse way of derivation, i.e., (3.77) from (3.78). For $r = 1$, by using (3.80), it can be shown that

$$\frac{p(Y_1 = 0, H_0) \cdot p(Y_1 = 1, H_0)}{p(Y_1 = 0, H_1) \cdot p(Y_1 = 1, H_1)} = 1. \quad (3.82)$$

For $r > 1$, from (3.81),

$$p(Y_1^{r+1} = y_1^r 0, H_0) \cdot p(Y_1^{r+1} = y_1^r 1, H_0) = p(Y_1^{r+1} = y_1^r 0, H_1) \cdot p(Y_1^{r+1} = y_1^r 1, H_1). \quad (3.83)$$

Therefore, (3.77) is equivalent to (3.78) for every $r \in \{1, \dots, N\}$.

We will show that if (3.79) can be satisfied, it also meets the the necessary and sufficient conditions (3.29)-(3.31) to achieve the Helstrom bound. Since (3.79) is derived from (3.78), which is equivalent to (3.77), it implies that if (3.77) is true, then the necessary and sufficient conditions (3.29)-(3.31) are satisfied.

From (3.79), it can be directly shown that

$$\begin{aligned} \frac{p(Y_1^N = y_1^{N-1} 0, H_0)}{p(Y_1^N = y_1^{N-1} 0, H_1)} &= \frac{1 - 2p_0 p_1 c^{2N} + \sqrt{1 - 4p_0 p_1 c^{2N}}}{2p_0 p_1 c^{2N}} \\ \frac{p(Y_1^N = z_1^{N-1} 1, H_0)}{p(Y_1^N = z_1^{N-1} 1, H_1)} &= \frac{1 - 2p_0 p_1 c^{2N} - \sqrt{1 - 4p_0 p_1 c^{2N}}}{2p_0 p_1 c^{2N}}, \end{aligned} \quad (3.84)$$

which also implies

$$\begin{aligned} \frac{p(Y_N = 0, H_0)}{p(Y_N = 0, H_1)} &= \frac{1 - 2p_0 p_1 c^{2N} + \sqrt{1 - 4p_0 p_1 c^{2N}}}{2p_0 p_1 c^{2N}} \\ \frac{p(Y_N = 1, H_0)}{p(Y_N = 1, H_1)} &= \frac{1 - 2p_0 p_1 c^{2N} - \sqrt{1 - 4p_0 p_1 c^{2N}}}{2p_0 p_1 c^{2N}}. \end{aligned} \quad (3.85)$$

By multiplying these two probabilities, (3.31), which is

$$\frac{P(Y_N = 0, H_0) \cdot P(Y_N = 1, H_0)}{P(Y_N = 0, H_1) \cdot P(Y_N = 1, H_1)} = 1$$

can be first shown.

From $p(Y_N = 0, H_0) + p(Y_N = 1, H_0) = p_0$ and $p(Y_N = 0, H_1) + p(Y_N = 1, H_1) = p_1$, the second equation of (3.85) can be written as

$$\frac{p_0 - p(Y_N = 0, H_0)}{p_1 - p(Y_N = 0, H_1)} = \frac{1 - 2p_0p_1c^{2N} - \sqrt{1 - 4p_0p_1c^{2N}}}{2p_0p_1c^{2N}}, \quad (3.86)$$

which is equivalent to

$$\begin{aligned} & 2p_0p_1c^{2N} (p_0 - p(Y_N = 0, H_0)) \\ &= (p_1 - p(Y_N = 0, H_1)) \left(1 - 2p_0p_1c^{2N} - \sqrt{1 - 4p_0p_1c^{2N}} \right). \end{aligned} \quad (3.87)$$

By using this and the first equation of (3.85), we get

$$\begin{aligned} p(Y_N = 0, H_0) &= \frac{p_0}{2} \left(1 + \frac{1 - 2p_1c^{2N}}{\sqrt{1 - 4p_0p_1c^{2N}}} \right), \\ p(Y_N = 1, H_0) &= \frac{p_0}{2} \left(1 - \frac{1 - 2p_1c^{2N}}{\sqrt{1 - 4p_0p_1c^{2N}}} \right). \end{aligned} \quad (3.88)$$

The ratio between these two probabilities becomes

$$\begin{aligned} \frac{p(Y_N = 1, H_0)}{p(Y_N = 0, H_0)} &= \frac{\sqrt{1 - 4p_0p_1c^{2N}} - (1 - 2p_1c^{2N})}{\sqrt{1 - 4p_0p_1c^{2N}} + (1 - 2p_1c^{2N})} \\ &= \frac{1 - 2p_1c^{2N}(1 + p_0 - p_1c^{2N}) - (1 - 2p_1c^{2N})\sqrt{1 - 4p_0p_1c^{2N}}}{2p_1^2c^{2N}(1 - c^{2N})}. \end{aligned} \quad (3.89)$$

From (3.84) and (3.89), we can show that (3.29) and (3.30) are satisfied., i.e.,

$$\frac{p_0}{p_1} \cdot \frac{P(Y_1^N = y_1^{N-1}0, H_1)}{P(Y_1^N = y_1^{N-1}0, H_0)} = \left(c^N - \sqrt{(1 - c^{2N}) \frac{P(Y_N = 1, H_0)}{P(Y_N = 0, H_0)}} \right)^2,$$

$$\frac{p_0}{p_1} \cdot \frac{P(Y_1^N = y_1^{N-1} \mathbf{1}, H_1)}{P(Y_1^N = y_1^{N-1} \mathbf{1}, H_0)} = \left(c^N + \sqrt{(1 - c^{2N}) \frac{P(Y_N = 0, H_0)}{P(Y_N = 1, H_0)}} \right)^2.$$

To show the equality, we use the following.

$$\begin{aligned} & 2c^{2N} - \frac{1}{p_1} + \frac{1}{p_1} \sqrt{1 - 4p_0 p_1 c^{2N}} \\ &= \frac{2}{p_1} \sqrt{\frac{1 - 2p_0 p_1 c^{2N} + 2p_1^2 c^{4N} - 2p_1 c^{2N} - (1 - 2p_1 c^{2N}) \sqrt{1 - 4p_0 p_1 c^{2N}}}{2}}, \\ & - \left(2c^{2N} - \frac{1}{p_1} - \frac{1}{p_1} \sqrt{1 - 4p_0 p_1 c^{2N}} \right) \\ &= \frac{2}{p_1} \sqrt{\frac{1 - 2p_0 p_1 c^{2N} + 2p_1^2 c^{4N} - 2p_1 c^{2N} + (1 - 2p_1 c^{2N}) \sqrt{1 - 4p_0 p_1 c^{2N}}}{2}}, \end{aligned} \tag{3.90}$$

which is true since the left-hand side terms are positive (because $\sqrt{1 - 4p_0 p_1 c^{2N}} > 1 - 2p_1 c^{2N} > 0$), and the squared of left-hand sides are equal to the squared of the right-hand side terms, respectively.

Therefore, the adaptive measurement that satisfies (3.77) also meet the necessary and sufficient conditions of the Helstrom bound.

Capacity of Coherent Detection

■ 4.1 Introduction: Detection of Optical Signals

We start by describing the optical channel of interest. Over a given period of time $t \in [0, T)$, we first consider a constant input to the channel, which is a *coherent state*, denoted by $|S\rangle$, where $S \in \mathbb{C}$. Here, coherent state can be understood as simply the light generated from a laser. In a noise-free environment, if one uses a photon counter to receive this optical signal, the output of the photon counter is a Poisson process, with rate $\lambda = |S|^2$, indicating the arrivals of individual photons. Clearly, one can generalize from a constant input to have $S(t)$, $t \in [0, T)$, which results in a non-homogeneous Poisson process at the output. The cost of transmitting such optical signals is naturally the average number of photons, which is equal to $\int_0^T |S(t)|^2 dt$. Here, without loss of generality, we set the scaling factors on the rate and photon counts to 1, ignoring issues with linear attenuation and efficiency of optical devices. Such receivers based on photon counters that detect the intensity of the optical signals are called direct detection receivers, and the resulting communication channel is called a Poisson channel. The capacity of the Poisson channel is well studied [38, 49].

Since a coherent state optical signal can be described by a complex number S , it is of interest to design coherent receivers that measure the phase of S , and thus allow information to be modulated on the phase. The following architecture, proposed by Kennedy, is a particular front end of the receiver, the output of which depends on the phase of S .

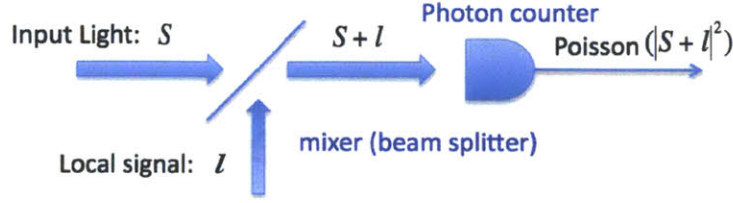


Figure 4.1. Coherent Receiver Using Local Feedback Signal

In Figure 4.1, instead of directly feeding the input optical signal $|S\rangle$ to the photon counter, a local signal $|l\rangle$ is mixed with the input, to generate a coherent state $|S + l\rangle$, and the output of the photon counter is a Poisson process with rate $|S + l|^2$. Note that l can in principle be chosen as an arbitrary complex number, with any desired phase difference from the input signal S . Thus, the output of this processing can be used to extract the phase information in the input. In a sense, the local signal is designed to control the channel through which the optical signal $|S\rangle$ is observed.

Kennedy used this receiver architecture to distinguish between binary hypotheses, i.e., two possible coherent states corresponding to waveforms $S_0(t), S_1(t), t \in [0, T)$, with prior probabilities π_0, π_1 , respectively, using a constant control signal l . This work was later generalized by Dolinar [10], where a control waveform $l(t), t \in [0, T)$ was used. The waveform $l(\cdot)$ is chosen adaptively based on the photon arrivals at the output. It was shown that the resulting probability of error for binary hypothesis testing is

$$P_e = \frac{1}{2} \left(1 - \sqrt{1 - 4\pi_0\pi_1 e^{-\int_0^T |S_0(t) - S_1(t)|^2 dt}} \right). \quad (4.1)$$

Somewhat surprisingly, this error probability coincides with the lower bound optimized over all possible quantum detectors, the Helstrom limit [21]. The optimality of Dolinar's receiver is an amazing result, as it shows that the minimum probability of error quantum detector for the binary problem can indeed be implemented with the very

simple receiver structure in Figure 4.1. Unfortunately, this result does not generalize to problems with more than two hypotheses.

The goal of the current chapter is twofold. We are interested in finding a natural generalization of Dolinar's receiver to general hypothesis testing problems with more than two possible signals. In addition, we would like to consider using such receivers to receive coded transmissions, and thus compute the information rate that can be reliably carried through the optical channel, with the above specific structure of the receiver front end. In the following, we will start by re-deriving Dolinar's design of the control waveform $l(t)$ to motivate our approach.

■ 4.2 Binary Hypothesis Testing

We consider the binary hypothesis testing problem with two possible input signals, $|S_0(t)\rangle, |S_1(t)\rangle$, under hypotheses $H = 0, 1$ respectively, and denote $\pi_0(t)$ and $\pi_1(t)$ as the posterior distribution over the two hypotheses, conditioned on the output of the photon counter up to time t . For simplicity, we assume that $S_0, S_1 \in \mathbb{R}$. Generalization to the complex-valued case will be straightforward. Based on the receiver knowledge, we choose the control signal $l(t)$, to be applied in an arbitrarily short interval $[t, t + \Delta)$. After observing the output during this interval, the receiver can update the posterior probabilities to obtain $\pi_0(t + \Delta)$ and $\pi_1(t + \Delta)$, and then follow the same procedure to choose the control signal in the next interval, and so on. As we pick Δ to be arbitrarily small, we can restrict the control signal $l(t)$ in such a short interval to be a constant l . In the following, we focus on solving the single step optimization of l in the above recursion, and drop the dependence on t to simplify the notation.

We first observe that the optimal value of l must be real, as having a non-zero imaginary part in l simply adds a constant rate to the two Poisson processes corresponding to the two hypotheses, and does not improve the quality of observation. We write $\lambda_i = (S_i + l)^2, i = 0, 1$ to denote the rate of the resulting Poisson processes. Over a very short period of time, the realized Poisson processes can have, with a high probability,

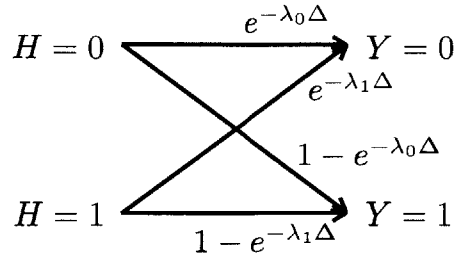


Figure 4.2. Effective binary channel between the input hypotheses and the observation over a Δ period of time

either 0 or 1 arrival, with probabilities $e^{-\lambda_i \Delta}$, $1 - e^{-\lambda_i \Delta}$, respectively.¹ Over this short period of time, the receiver front end can be thought of as a binary channel as shown in Figure 4.2. Note that the channel parameters λ_i 's depend on the value of the control signal l . Our goal is to pick an l for each short interval such that they contribute to the overall decision in the best way.

The difficulty here is that it is not obvious how we should quantify the “contribution” of the observation over a short period of time to the overall decision making. An intuitive approach one can use is to choose l that maximizes the mutual information over the binary channel. For convenience, we write the input to the channel as H and the output of the channel as $Y \in \{0, 1\}$, indicating either 0 or 1 photon arrival. The following result gives the solution to this optimization problem.

Lemma 4.1. *The optimal choice that maximizes the mutual information $I(H; Y)$ for*

¹One has to be careful in using the above approximation of the binary channel. As we are optimizing over the control signal, it is not obvious that the resulting λ_i 's are bounded. In other words, the mean of the Poisson distributions, $\lambda_i \Delta$, might not be small. The assumption of either 0 or 1 arrival, and the approximation in the corresponding probabilities, can be justified in two senses: First, a practical photon detector can easily sense whether or not any number of photons arrives during a time interval Δ , but cannot exactly count the number of photon arrivals, especially as $\Delta \rightarrow 0$. So, the binary output channel model is much more practical than the Poisson output channel model. Second, when we want to maximize the ability to distinguish between two hypotheses $H = 0, 1$, we essentially need to distinguish between the signal amplitudes S_0 and S_1 through photon arrivals. Adding the feedback control signal $l \rightarrow \infty$ does not help in distinguishing S_0 and S_1 . In this sense, we can reason that the optimum l would not make λ_i unbounded.

the effective binary channel is

$$l^* = \frac{S_0\pi_0 - S_1\pi_1}{\pi_1 - \pi_0}. \quad (4.2)$$

With this choice of the control signal, the following relation holds:

$$\pi_0\sqrt{\lambda_0} = \pi_1\sqrt{\lambda_1}. \quad (4.3)$$

Remark 4.1. Note that the optimal choice of l that maximizes the mutual information $I(H; Y)$ is exactly the same as the choice of l in (3.46) that is optimized with a totally different view of minimizing instant detection error probability. It turns out that for “binary” hypothesis testing, the feedback control signal l that maximizes mutual information of the corresponding channel coincides with that of minimizing detection error probability at each instant. This is not generally true for M -ary ($M > 2$) hypothesis testing.

Proof. Appendix 4.A □

The relation in (4.3) gives some useful insights. If $\pi_0 > \pi_1$, we have $\lambda_1 > \lambda_0$, and vice versa. That is, by switching the sign of the control signal l , we always make the Poisson rate corresponding to the hypothesis with the higher probability smaller. In the short interval where this control is applied, with a high probability we would observe no photon arrival, in which case we would confirm the more likely hypothesis. For a very small value of Δ , this occurs with a dominating probability, such that the posterior distribution moves only by a very small amount. On the other hand, when there is indeed an arrival, i.e., $Y = 1$, we would be quite surprised, and the posterior distribution of the hypotheses moves away from the prior. Considering this latter case, the updated distribution over the hypotheses can be written as

$$\frac{\Pr(H = 1|Y = 1)}{\Pr(H = 0|Y = 1)} = \frac{\pi_1 \cdot \lambda_1 \Delta}{\pi_0 \cdot \lambda_0 \Delta} + O(\Delta) = \frac{\pi_0}{\pi_1} + O(\Delta).$$

The posterior distribution under the case of 0 or 1 arrival turns out to be inverse to each other as $\Delta \rightarrow 0$. In other words, the larger one of the two probabilities of the hypotheses remains the same no matter if there is an arrival in the interval or not. As we apply such optimal control signals recursively, this larger value progresses towards 1 at a predictable rate, regardless of when and how many arrivals are observed. *The random photon arrivals only affect the decision on which is the more likely hypothesis, but do not affect the quality of this decision.* The next lemma describes this recursive control signal and the resulting performance. Without loss of generality, we assume that at $t = 0$, the prior distribution satisfies $\pi_0 \geq \pi_1$. Also we let $N(t)$ be the number of arrivals observed in $[0, t)$.

Lemma 4.2. *Let $g(t)$ satisfy $g(0) = \pi_0/\pi_1$ and*

$$g(t) = g(0) \cdot \exp \left[\int_0^t \frac{(S_0(\tau) - S_1(\tau))^2 (g(\tau) + 1)}{g(\tau) - 1} d\tau \right].$$

The recursive mutual-information maximization procedure described above yields a control signal

$$l^*(t) = \begin{cases} l_0(t) & \text{if } N(t) \text{ is even} \\ l_1(t) & \text{if } N(t) \text{ is odd} \end{cases}$$

where

$$l_0(t) = \frac{S_1(t) - S_0(t)g(t)}{g(t) - 1}, \quad l_1(t) = \frac{S_0(t) - S_1(t)g(t)}{g(t) - 1}.$$

Furthermore, at time T , the decision of the hypothesis testing problem is $\hat{H} = 0$ if $N(T)$ is even, and $\hat{H} = 1$ otherwise. The resulting probability of error coincides with (4.1).

Proof. Appendix 4.B □

Figure 4.3 shows an example of the optimal control signal. The plot is for a case where $S_i(t)$'s are constant on-off-keying waveforms. As shown in the plot, the control

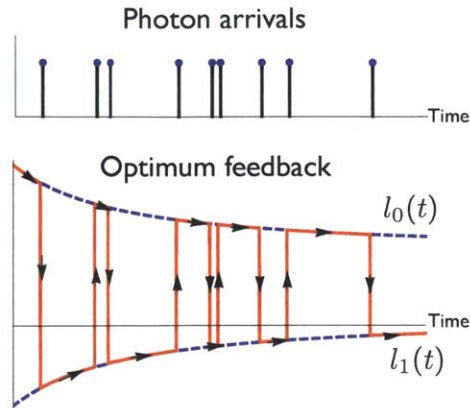


Figure 4.3. An example of the control signal that achieves the minimum probability of error.

signal $l(t)$ jumps between two prescribed curves, $l_0(t), l_1(t)$, corresponding to the cases $\pi_0(t) > \pi_1(t)$ and $\pi_0(t) < \pi_1(t)$, respectively. With the proper choice of the control signal, each time when there is a photon arrival, the receiver is so surprised that it flips its choice of \hat{H} . However, $g(t) = \max\{\pi_0(t), \pi_1(t)\} / \min\{\pi_0(t), \pi_1(t)\}$, indicating how much the receiver is committed to the more likely hypothesis, increases at a prescribed rate regardless of the arrivals.

■ 4.3 Generalization to M-ary Hypothesis Testing

The success in the binary hypothesis testing problem reveals some useful insights for general dynamic communication problems. Regardless of the physical channel that one communicates over, one can always have a “slow motion” understanding of the process by studying how the posterior distribution over the messages evolves over time. Over the process of communications, this posterior distribution, conditioned on more and more observations of the receiver, should move from the prior towards a deterministic distribution, allowing the receiver to “lock in” on a particular message. This viewpoint is more general than the conventional setup in information theory, and particularly useful in understanding dynamic problems, as it is not based on any notion of sufficient statistics, block codes, or any predefined notion of reliability. As we measure how far

the posterior distribution moves at each time, we can quantify how the communication process at each time point contributes to the overall decision making.

The optimality result in Lemma 4.2 is, however, difficult to duplicate for general M -ary problems. Of course, we can always mimic the procedure, to choose the control signal that maximizes the mutual information over an M -input-binary-output channel at each time. The result does not always give the minimum probability of error in general. The reason for that is intuitive. There is a fundamental difference between maximizing mutual information and minimizing the probability of error. In other words, on a general M -ary alphabet, a posterior distribution with a lower entropy does not necessarily have a lower probability of detection errors. These two coincide only for the binary case, since the posterior distribution over two messages lives in a single-dimensional space. In general, the goal of decision making favors the posterior distributions, over the messages, with a dominating largest element; maximizing mutual information, however, does not distinguish what kind of information is conveyed.

Consequently, it is hard to define a metric on the efficiency of communication over a time interval in the middle of a communication session that precisely measures how well this interval serves the overall purpose. Even if one can define a precise metric, it is often hard to imagine that the analytical solution of the optimal control signal or the resulting performance can be computed from optimizing such metrics. Moreover, such metrics should be time-varying, depending on how much time is left before the decision is made. Intuitively, at an early time point (i.e., when a longer time remains before a final decision needs to be made), since the current observation will be combined with a lot more information yet to come, we are more willing to take any kind of information, and hence it makes sense to maximize mutual information. On the other hand, as the deadline of decision making approaches, the system becomes much “pickier” in choosing information to extract from measurements, and demands only the information that helps the receiver to lock in to one particular message. Thus, the control signal should be optimized accordingly.

To test this intuition, we restrict our attention to the family of Rényi entropy. Rényi entropy of order α of a given distribution P over an alphabet \mathcal{X} is defined as

$$H_\alpha(P) = \frac{1}{1-\alpha} \log \left(\sum_{x \in \mathcal{X}} P^\alpha(x) \right).$$

It is easy to verify that as $\alpha \rightarrow 1$, $H_\alpha(P)$ is the Shannon entropy, and as $\alpha \rightarrow \infty$, $H_\infty(P) = -\log(\max_{x \in \mathcal{X}} P(x))$, which is a measure of the probability of error in guessing the value of X , with distribution P , as $\hat{X} = \arg \max_x P(x)$.

Now for general M -ary hypothesis testing problems, we consider a recursive design of the control signal l similar to that introduced in section 4.2, except that at each time, rather than maximizing the mutual information over the effective channel, which is equivalent to minimizing the conditional Shannon entropy of the messages, we instead minimize the average Rényi α entropy, i.e., we solve the optimization problem

$$\min_l \sum_y P_Y(y) \cdot H_\alpha(P_{H|Y=y}(\cdot)). \quad (4.4)$$

Intuitively, for $\alpha \in [1, \infty)$, the larger α is, optimization in (4.4) tends more in favor of posterior distributions that are concentrated on a single entry. Smaller values of α , on the other hand, correspond to receivers that are more agnostic to what type of information is obtained. A good design should use smaller values of α at the beginning of the communication session (when a longer time remains before the final decision is made), and increase α as a shorter time remains before the final decision. In the numerical example in Figure 4.4, to illustrate the point, we compare the cases where α is chosen to be fixed to be either 1 or 100 throughout the time $t \in [0, T]$. Our intuition says that choosing a smaller α is desirable, when we have enough time to collect information before the final decision. On the other hand, when we need to make a final decision immediately, a larger α is preferable. We observe that using $\alpha = 1$ gives better performance if T is longer, and choosing a larger α yields better error probability when T is small. This experiment confirms our intuition.

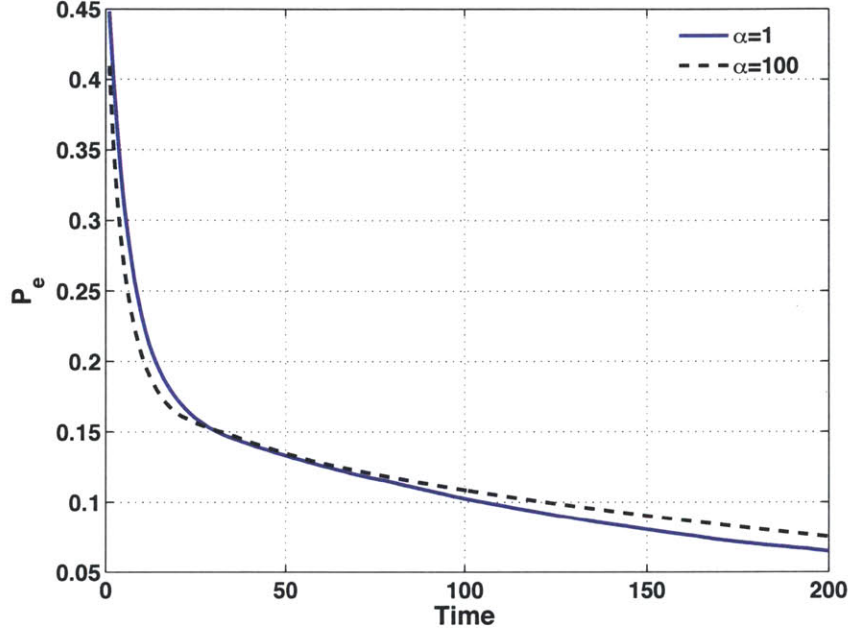


Figure 4.4. Empirical average of detection error probability (after 10,000 runs) for ternary hypothesis testing, using control signals that minimize the average Rényi α -entropy for different values of α ; Ternary inputs $\{|5\rangle, |-6\rangle, |3\rangle\}$ are used with prior probabilities $p = \{0.5, 0.4, 0.1\}$.

■ 4.4 Coded Transmissions and Capacity Results

We now turn our attention to the problem of coded transmissions over the optical channel with coherent receivers. We are interested in finding the classical capacity of such channels, i.e., the number of information bits that can be modulated into the optical signals, and reliably decoded with the receiver architecture shown in Figure 4.1. We are particularly interested in the case where the average number of transmitted photons is small, and hence a high photon efficiency, in bits/photon, is achieved.

The capacity of the same channel without the constraint in the receiver architecture is studied in [23, 37]. It is shown [14] that the capacity of the channel is given by

$$C_{\text{Holevo}}(\mathcal{E}) = (1 + \mathcal{E}) \log(1 + \mathcal{E}) - \mathcal{E} \log \mathcal{E} \text{ nats/channel use} \quad (4.5)$$

where \mathcal{E} is the average number of photons transmitted per channel use. To achieve this

data rate, an optimal joint quantum measurement over a long sequence of symbols must be used. In practice, however, such measurement is very hard to implement. We are therefore interested in finding the achievable data rate when a simple receiver structure is adopted. Nevertheless, (4.5) serves as a performance benchmark. In the regime of interests where $\mathcal{E} \rightarrow 0$, it is useful to approximate (4.5) as

$$C_{\text{Holevo}}(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} + \mathcal{E} + o(\mathcal{E}). \quad (4.6)$$

As another performance benchmark, we also consider the capacity when a direct detection receiver is used. The capacity of this channel is studied in [38, 49], and the regime of low average photon numbers is studied in [29]. For our purpose of performance comparison, we actually need a more precise scaling law of performance. The following lemma describes such a result.

Lemma 4.3 (Capacity of Direct Detection). *As $\mathcal{E} \rightarrow 0$, the optimal input distribution to the optical channel with direct detection is on-off-keying, with*

$$|S\rangle = \begin{cases} |0\rangle, & \text{with prob. } 1 - p^* \\ |\sqrt{\mathcal{E}/p^*}\rangle, & \text{with prob. } p^* \end{cases}$$

where $\lim_{\mathcal{E} \rightarrow 0} \frac{p^*}{\frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}} = 1$, and the resulting capacity is

$$C_{\text{DD}}(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}). \quad (4.7)$$

Proof. Appendix 4.C. □

Comparing (4.6) and (4.7), we observe that the two capacities have the same leading term. This means as $\mathcal{E} \rightarrow 0$, the optimal photon efficiency of $\log(1/\mathcal{E})$ nats/photon can be achieved even with a very simple direct detection receiver.

In practice, however, the two performances have significant differences. For example, if one wishes to achieve a photon efficiency of 10 bits/photon, one can solve for \mathcal{E}

that satisfies $C(\mathcal{E})/\mathcal{E} = 10$ bits/photon in both cases, and get $\mathcal{E}_{\text{Holevo}} \approx 0.0027$ and $\mathcal{E}_{\text{DD}} \approx 0.00010$. The resulting capacities (bits/channel use), or (equivalently) spectral efficiencies, also differ by more than one order of magnitude. This example indicates that although (4.6) and (4.7) have the same limit as $\mathcal{E} \rightarrow 0$, the rates at which this limit is approached are quite different, which is of practical importance. Similar phenomena have also been observed for wideband wireless channels [45, 52].

As a result, the 2nd term in the capacity results cannot be ignored. In fact, any reasonable scheme with coherent processing should at least achieve a rate higher than that with direct detection, and thus should have the leading term as $\mathcal{E} \log \frac{1}{\mathcal{E}}$. It is the second term in the achievable rate that indicates whether a new scheme is making a significant step towards achieving the Holevo capacity limit. In the following, we will study the achievable rates over the optical channel with receiver front end as shown in Figure 4.1, and evaluate the performance according to this scaling law.

The problem of coded transmission and finding the maximum information rate that can be conveyed through an optical channel with a coherent processing receiver is in fact easier than that of hypothesis testing, even though there are exponentially many possible messages. One first observation is that when communicating with a long block of N symbols, there is no issue of a pressing deadline of decision for most of the time. Therefore, it makes sense to always use the mutual information maximization to decide which control signal to apply. A straightforward generalization of the Dolinar receiver can be described as follows:

First, at each time instance $i \in \{1, \dots, N\}$, the encoding map can be written $f_i : \{1, 2, \dots, M = 2^{NR}\} \rightarrow X_i \in \mathcal{X}$, where X_i is the symbol transmitted in the i^{th} use of the channel. This map ensures that X_i has a desired input distribution P_X , computed under the assumption that all messages are equally likely, i.e., $\frac{1}{2^{NR}} |\{m : f_i(m) = x\}| = P_X(x)$, $\forall x \in \mathcal{X}$.

The receiver keeps track of the posterior distribution over the messages. Given $P_{M|Y^{i-1}}(\cdot|y^{i-1})$, which is the distribution over the messages conditioned on the previous

observations, the effective input distribution $P'_X(x) = \sum_{m: f_i(m)=x} P_{M|Y^{i-1}}(m|y^{i-1})$ can be computed. Using this as the prior distribution of the transmitted symbol, one can apply the control signal that maximizes the mutual information.

Upon observing the output Poisson process in the i^{th} symbol period, denoted as $Y_i = y_i$, the receiver computes the posterior distribution of the transmitted symbol $P''_X(x) = P_{X_i|Y_i}(x|y_i)^2$, and uses that to update its knowledge of the messages:

$$P_{M|Y^i}(m|y^i) = P_{M|Y^{i-1}}(m|y^{i-1}) \cdot \frac{P''_X(x)}{P'_X(x)},$$

for all m such that $f_i(m) = x$.

Repeating this process, we have a coherent-processing receiver based on updating the receiver knowledge. There are two further simplifications that make the analysis of this scheme even simpler.

First, we observe that with exponentially many messages, for a dominating fraction of the time when the block code is transmitted, the receiver's knowledge, $P_{M|Y^i}$, is such that the probability of any message, including the correct one, is exponentially small. Thus, with a random coding map f_i , P'_X is very close to P_X . Thus, the step of updating the receiver's knowledge is in fact not important. This assumption starts to fail only when the receiver starts to lock in a specific message, i.e., when $P_{M|Y^i}(m)$ is not exponentially small for some m . It is shown in [7] that the fraction of time when this happens is indeed very small, and can thus be ignored when a long-term average performance metric such as the data rate is of concern.

Second, suppose we choose the optimal input distribution, which maximizes the photon efficiency, over a short period Δ of time. After using this input for Δ time, the receiver would update the posterior distribution, which makes the effective input distribution on X deviate from the optimum. This is undesirable. One can avoid this problem by using very short symbol periods. That is, after transmitting for a very short

²We omit the conditioning on the history Y^{i-1} here to emphasize that the update is based on the observations in a single symbol period.

time period, the transmitter should re-shuffle the messages so that the distribution of the transmitted symbols, conditioned on the receiver's knowledge, is re-adjusted back to the optimal choice. This is precisely the same argument we used regarding classical communication over wide-band channels. As a result, we do not have to worry about updating the receiver's knowledge and the control signals even within a symbol period. Instead, we are interested only in the photon efficiency over a short time period. In other words, we can focus only on a thin slice on the left end of Figure 4.4.

Based on these observations, we state our results in the photon efficiency of the optical channel of interest.

Theorem 4.4. *For the optical channel with a receiver front end as shown in Figure 4.1, and sequentially updated control signals, suppose that the transmitted symbols are drawn from a finite alphabet, i.e., at each time the transmitted optical signal $|X_i\rangle$ is chosen from $X_i \in \mathcal{X} \subset \mathbb{C}$, with $|\mathcal{X}|$ finite. Then the achieved photon efficiency is upper bounded by*

$$\frac{C(\mathcal{E})}{\mathcal{E}} \leq \log \frac{1}{\mathcal{E}} - \log \log \frac{1}{\mathcal{E}} + O(1). \quad (4.8)$$

Proof. Appendix 4.D. □

This says that essentially the achievable photon efficiency with coherent receivers is not significantly different from that of direct detection receivers.

This theorem is a useful step in understanding more general coherent receivers, with joint processing of multiple symbols. Here, we describe a general optical receiver with *classical processing* as follows. Let the codeword transmitted by a sequence of coherent states $|X_1\rangle|X_2\rangle \dots |X_N\rangle$, where each X_i is drawn from a finite alphabet. The receiver forms $|Y_1\rangle, |Y_2\rangle, \dots, |Y_M\rangle$ and uses a photon counter to observe them separately. Each Y_j is formed by a linear passive mixing of the X_i 's and an arbitrary control signal l_j : $Y_j = \sum_{i=1}^N \alpha_{ij} X_i + l_j$, where α_{ij} satisfy $\sum_j |\alpha_{ij}|^2 \leq 1, \forall i$ and $\sum_i |\alpha_{ij}|^2 \leq 1, \forall j$, which ensure the physical constraint of energy conservation and the fact that duplication

or noiseless amplification of coherent states are not possible. The mixing parameters and the control signals can be decided sequentially based on the earlier observations. Following the spirit of Theorem 4.4, we state the following conjecture.

Conjecture 4.1. *The achievable photon efficiency by an optical receiver with classical processing satisfies (4.8).*

While this conjecture is negative by nature, it is of practical importance. It implies that in order to achieve the photon efficiency predicted by the Holevo limit, it is necessary to resort to quantum processing that introduces non-classical states, such as entangled or squeezed states. The approach of mixing coherent states and applying a local control signal would not yield significant improvement in photon efficiency.

■ 4.5 Conclusion

We studied the general coherent-state hypothesis testing problem and the capacity of the pure-loss optical channel with a *coherent processing* receiver (a receiver that uses coherent feedback control and direct detection). We described the Dolinar receiver with the intention of optimizing the communication efficiency at each instant, based on recursively updated knowledge of the receiver. Using this viewpoint, we gave a natural generalization of the design to general M -ary hypothesis testing problems. We analyzed the information capacity with coherent receivers, and compared the result with that of direct detection receivers and of arbitrary quantum receivers (the Holevo limit), using the appropriate scalings in the low photon number regime.

While Theorem 4.4, is a negative result by its nature, it is of practical importance. It implies that in order to achieve the photon efficiency predicted by the Holevo limit, it is necessary to resort to complicated quantum processing that may include entanglement or squeezing. The approach of mixing coherent states and applying a local control signal would not yield significant improvement in terms of photon efficiency. However, the proposed sequential receiver designs and the adaptive feedback control of the receiver can be applied for more general dynamic communication problems.

■ 4.A Proof of Lemma 4.1

For input distribution $\{\pi_0, \pi_1\}$, the mutual information between the binary input H and the binary output Y is

$$\begin{aligned}
I(H; Y) &= \pi_0 \left(e^{-\lambda_0 \Delta} \cdot \log \frac{e^{-\lambda_0 \Delta}}{\pi_0 e^{-\lambda_0 \Delta} + \pi_1 e^{-\lambda_1 \Delta}} + (1 - e^{-\lambda_0 \Delta}) \log \frac{1 - e^{-\lambda_0 \Delta}}{1 - \pi_0 e^{-\lambda_0 \Delta} - \pi_1 e^{-\lambda_1 \Delta}} \right) \\
&\quad + \pi_1 \left(e^{-\lambda_1 \Delta} \cdot \log \frac{e^{-\lambda_1 \Delta}}{\pi_0 e^{-\lambda_0 \Delta} + \pi_1 e^{-\lambda_1 \Delta}} + (1 - e^{-\lambda_1 \Delta}) \log \frac{1 - e^{-\lambda_1 \Delta}}{1 - \pi_0 e^{-\lambda_0 \Delta} - \pi_1 e^{-\lambda_1 \Delta}} \right).
\end{aligned} \tag{4.9}$$

As $\Delta \rightarrow 0$, it can be approximated as

$$I(H; Y) = (\pi_0 \lambda_0 \log \lambda_0 + \pi_1 \lambda_1 \log \lambda_1 - (\pi_0 \lambda_0 + \pi_1 \lambda_1) \log(\pi_0 \lambda_0 + \pi_1 \lambda_1)) \Delta + O(\Delta^2). \tag{4.10}$$

To find the optimum feedback control signal l that maximizes $I(H; Y)$, we take the derivative of $I(H; Y)$ over l .

$$\begin{aligned}
\frac{d}{dl} I(H; Y) &= \left(\frac{d}{d\lambda_0} I(H; Y) \right) \cdot \left(\frac{d}{dl} \lambda_0 \right) + \left(\frac{d}{d\lambda_1} I(H; Y) \right) \cdot \left(\frac{d}{dl} \lambda_1 \right) \\
&= (\pi_0 \log \lambda_0 - \pi_0 \log(\pi_0 \lambda_0 + \pi_1 \lambda_1)) 2\sqrt{\lambda_0} \Delta \\
&\quad + (\pi_1 \log \lambda_1 - \pi_1 \log(\pi_0 \lambda_0 + \pi_1 \lambda_1)) 2\sqrt{\lambda_1} \Delta + O(\Delta^2) \\
&= 2\Delta \left(\pi_0 \sqrt{\lambda_0} (\log \lambda_0 - \log(\pi_0 \lambda_0 + \pi_1 \lambda_1)) + \pi_1 \sqrt{\lambda_1} (\log \lambda_1 - \log(\pi_0 \lambda_0 + \pi_1 \lambda_1)) \right) \\
&\quad + O(\Delta^2)
\end{aligned} \tag{4.11}$$

Note that when

$$l = \frac{S_0 \pi_0 - S_1 \pi_1}{\pi_1 - \pi_0}, \tag{4.12}$$

and thus when $\pi_0\sqrt{\lambda_0} = \pi_1\sqrt{\lambda_1}$, the first term of $\frac{d}{dt}I(H;Y)$ becomes 0 since

$$\begin{aligned} & \pi_0\sqrt{\lambda_0}(\log \lambda_0 - \log(\pi_0\lambda_0 + \pi_1\lambda_1)) + \pi_1\sqrt{\lambda_1}(\log \lambda_1 - \log(\pi_0\lambda_0 + \pi_1\lambda_1)) \\ &= \pi_0\sqrt{\lambda_0} \cdot \log \frac{\lambda_0 \cdot \lambda_1}{(\pi_0\lambda_0 + \pi_1\lambda_1)^2} = \pi_0\sqrt{\lambda_0} \cdot \log \frac{1}{(\pi_1 + \pi_0)^2} = 0. \end{aligned} \quad (4.13)$$

Therefore, the feedback control signal l in (4.12) maximizes $I(H;Y)$.

■ 4.B Proof of Lemma 4.2

Without the loss of generality, we assume that $\pi_0 \geq \pi_1$, and let $g(0) = \pi_0/\pi_1$. To understand how the optimum $l(t)$ should be updated over time, let us first assume that $S_0(t)$, $S_1(t)$, and $l(t)$ do not change during every Δ -interval, i.e., for $t \in [k \cdot \Delta, (k+1) \cdot \Delta)$, $k \in \{0, 1, \dots\}$. We will make $\Delta \rightarrow 0$ afterward. Let $Y_k \in \{0, 1\}$ indicate whether or not photons arrive during $t \in [k \cdot \Delta, (k+1) \cdot \Delta)$. When $Y_0 = 0$,

$$\frac{p(H=0|Y_0=0)}{p(H=1|Y_0=0)} = \frac{\pi_0}{\pi_1} \cdot \frac{p(Y_0=0|H=0)}{p(Y_0=0|H=1)} = \frac{\pi_0}{\pi_1} \cdot \frac{e^{-(S_0(0)+l(0))^2\Delta}}{e^{-(S_1(0)+l(0))^2\Delta}}. \quad (4.14)$$

As shown in Lemma 4.1, the optimum $l(0)$ is

$$l(0) = \frac{S_0(0)\pi_0 - S_1(0)\pi_1}{\pi_1 - \pi_0}. \quad (4.15)$$

When we plug this value into (4.14),

$$\frac{p(H=0|Y_0=0)}{p(H=1|Y_0=0)} = \frac{\pi_0}{\pi_1} \cdot \exp \left[(S_0(0) - S_1(0))^2 \cdot \frac{g(0) + 1}{g(0) - 1} \cdot \Delta \right]. \quad (4.16)$$

For this case $p(H=0|Y_0=0)/p(H=1|Y_0=0) \geq 1$.

When photons arrive during the first Δ -interval, i.e., $Y_0 = 1$, the posterior distribution over the hypotheses can be written as

$$\frac{p(H=0|Y_0=1)}{p(H=1|Y_0=1)} = \frac{\pi_0}{\pi_1} \cdot \frac{p(Y_0=1|H=0)}{p(Y_0=1|H=1)} = \frac{\pi_0}{\pi_1} \cdot \frac{1 - e^{-(S_0(0)+l(0))^2\Delta}}{1 - e^{-(S_1(0)+l(0))^2\Delta}}. \quad (4.17)$$

As $\Delta \rightarrow 0$,

$$\frac{p(H = 0|Y_0 = 1)}{p(H = 1|Y_0 = 1)} = \frac{\pi_0}{\pi_1} \cdot \frac{(S_0(0) + l(0))^2}{(S_1(0) + l(0))^2} + O(\Delta) = \frac{\pi_1}{\pi_0} + O(\Delta). \quad (4.18)$$

The posterior distribution under $Y_0 = 0$ in (4.16) and $Y_0 = 1$ in (4.18) turns out to be approximately the inverse of each other as $\Delta \rightarrow 0$. Therefore, when we denote the posterior distributions of $H = 0$ and $H = 1$ at time t as $\pi_0(t)$ and $\pi_1(t)$, respectively,

$$g(t) := \max\{\pi_0(t)/\pi_1(t), \pi_1(t)/\pi_0(t)\} \quad (4.19)$$

is uniquely determined regardless of the photon arrival history during $[0, t)$. Thus, for ease of calculation of $g(t)$, we will assume no photon arrivals during $[0, t)$, and evaluate the ratio between the posterior distributions of the two hypotheses.

First of all, from (4.16),

$$g(\Delta) = \frac{\pi_0}{\pi_1} \cdot \exp \left[(S_0(0) - S_1(0))^2 \cdot \frac{g(0) + 1}{g(0) - 1} \cdot \Delta \right]. \quad (4.20)$$

If we assume no photon arrival for the next $(N - 1)$ -intervals, i.e., $Y_1 = \dots = Y_{N-1} = 0$, then $g(k \cdot \Delta) = p(H = 0|Y_0^{k-1} = \mathbf{0})/p(H = 1|Y_0^{k-1} = \mathbf{0})$, and we get the following recursive equation for $g(N \cdot \Delta)$,

$$g(N \cdot \Delta) = \frac{\pi_0}{\pi_1} \cdot \exp \left[\sum_{k=0}^{N-1} \left((S_0(k \cdot \Delta) - S_1(k \cdot \Delta))^2 \cdot \frac{g(k \cdot \Delta) + 1}{g(k \cdot \Delta) - 1} \cdot \Delta \right) \right]. \quad (4.21)$$

By taking $\Delta \rightarrow 0$,

$$\begin{aligned} g(t) &= \frac{\pi_0}{\pi_1} \cdot \exp \left[\int_0^t \left((S_0(\tau) - S_1(\tau))^2 \cdot \frac{g(\tau) + 1}{g(\tau) - 1} \right) d\tau \right] \\ &= g(0) \cdot \exp \left[\int_0^t \left((S_0(\tau) - S_1(\tau))^2 \cdot \frac{g(\tau) + 1}{g(\tau) - 1} \right) d\tau \right]. \end{aligned} \quad (4.22)$$

The optimum feedback control signal at time t is equal to

$$l^*(t) = \frac{S_0\pi_0(t) - S_1\pi_1(t)}{\pi_1(t) - \pi_0(t)}. \quad (4.23)$$

Note that whenever a photon arrives, $\pi_0(t)/\pi_1(t)$ gets flipped. Therefore, starting from $\pi_0 \geq \pi_1$, when the number of photon arrivals until time t , which is denoted as $N(t)$, is an even integer, then $\pi_0(t) \geq \pi_1(t)$ so that $g(t) = \pi_0(t)/\pi_1(t)$. If $N(t)$ is odd, then $\pi_1(t) \geq \pi_0(t)$, which results in $g(t) = \pi_1(t)/\pi_0(t)$. Therefore,

$$l^*(t) = \begin{cases} l_0(t) & \text{if } N(t) \text{ is even} \\ l_1(t) & \text{if } N(t) \text{ is odd} \end{cases}$$

where

$$l_0(t) = \frac{S_1(t) - S_0(t)g(t)}{g(t) - 1}, \quad l_1(t) = \frac{S_0(t) - S_1(t)g(t)}{g(t) - 1}.$$

Furthermore, the decision of the hypothesis testing problem is $\hat{H} = 0$ if $N(T)$ is even, and $\hat{H} = 1$ otherwise. The average probability of error is then equal to $P_e = \min\{\pi_0(t), \pi_1(t)\}$, and by the definition of $g(t)$,

$$P_e = \frac{1}{1 + g(t)}. \quad (4.24)$$

It can also be shown that

$$g(t) = \frac{(1 + g(0))^2}{2g(0)} e^{m(t)} - 1 + \frac{1 + g(0)}{2g(0)} \sqrt{(1 + g(0))^2 e^{2m(t)} - 4g(0)e^{m(t)}} \quad (4.25)$$

with $m(t) = \int_0^t (S_0(\tau) - S_1(\tau))^2 d\tau$ is the solution for $g(t)$ in (4.22). And the resulting P_e is

$$P_e = \frac{1}{1 + g(t)} = \frac{1}{2} \left[1 - \sqrt{1 - 4\pi_0\pi_1 e^{-\int_0^t (S_0(\tau) - S_1(\tau))^2 d\tau}} \right]. \quad (4.26)$$

■ 4.C Proof of Lemma 4.3

The converse part of this lemma, i.e., that any finite alphabet input distribution to the optical channel with direct detection can never achieve information rates greater than

$$\mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}), \quad (4.27)$$

can be inferred from the converse proof of Theorem 4.4, which considers a more general receiver type, which makes the direct detection just a special case. Therefore, here we will only prove the achievability of the capacity in (4.27) with on-off-keying.

With the on-off keying input

$$|S\rangle = \begin{cases} |0\rangle, & \text{with prob. } 1-p \\ |\sqrt{\mathcal{E}/p}\rangle, & \text{with prob. } p, \end{cases}$$

the mutual information of the binary input and output channel generated by direct detection can be written as

$$\begin{aligned} & H_{\mathbf{B}} \left(1 - p + e^{-\frac{\mathcal{E}}{p}} \cdot p \right) - p \cdot H_{\mathbf{B}} \left(e^{-\frac{\mathcal{E}}{p}} \right) \\ &= - \left(1 - p + e^{-\frac{\mathcal{E}}{p}} \cdot p \right) \log \left(1 - p + e^{-\frac{\mathcal{E}}{p}} \cdot p \right) - p \left(1 - e^{-\frac{\mathcal{E}}{p}} \right) \log \left(p \left(1 - e^{-\frac{\mathcal{E}}{p}} \right) \right) \\ & \quad + p \cdot e^{-\frac{\mathcal{E}}{p}} \log \left(p \cdot e^{-\frac{\mathcal{E}}{p}} \right) + p \left(1 - e^{-\frac{\mathcal{E}}{p}} \right) \log \left(1 - e^{-\frac{\mathcal{E}}{p}} \right) \end{aligned} \quad (4.28)$$

where $H_{\mathbf{B}}(p) = -p \log p - (1-p) \log(1-p)$.

For $\lim_{\mathcal{E} \rightarrow 0} \frac{p}{\frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}} = 1$, we can approximate

$$\begin{aligned} e^{-\frac{\mathcal{E}}{p}} &= 1 - \frac{2}{\log(1/\mathcal{E})} + O\left(\frac{1}{(\log(1/\mathcal{E}))^2}\right), \\ 1 - p + e^{-\frac{\mathcal{E}}{p}} \cdot p &= 1 - \mathcal{E} + o(\mathcal{E}), \\ p \cdot e^{-\frac{\mathcal{E}}{p}} &= \mathcal{E} + O\left(\frac{\mathcal{E}}{\log(1/\mathcal{E})}\right). \end{aligned} \quad (4.29)$$

And by using these approximations, we can show that

$$\begin{aligned} H_{\mathbb{B}}\left(1 - p + e^{-\frac{\mathcal{E}}{p}} \cdot p\right) &= \mathcal{E} \log \frac{1}{\mathcal{E}} + O(\mathcal{E}), \\ p \cdot H_{\mathbb{B}}\left(e^{-\frac{\mathcal{E}}{p}}\right) &= \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}), \end{aligned} \quad (4.30)$$

and thus

$$H_{\mathbb{B}}\left(1 - p + e^{-\frac{\mathcal{E}}{p}} \cdot p\right) - p \cdot H_{\mathbb{B}}\left(e^{-\frac{\mathcal{E}}{p}}\right) = \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}). \quad (4.31)$$

Therefore, with on-off keying inputs and direction detection, we can achieve

$$C_{\text{DD}}(\mathcal{E}) \geq \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}). \quad (4.32)$$

■ 4.D Proof of Theorem 4.4

Let us first calculate the maximum mutual information of the binary input channel under the average photon number constraint of \mathcal{E} . We want to find the binary input $\{|S_0\rangle, |S_1\rangle\}$ with prior distribution $\{1 - p, p\}$ that satisfies $(1 - p)|S_0|^2 + p|S_1|^2 = \mathcal{E}$ while maximizing

$$\begin{aligned} \max_l I(H; Y) &= H_{\mathbb{B}}\left(1 - (1 - p)e^{-|S_0+l|^2} - pe^{-|S_1+l|^2}\right) \\ &\quad - (1 - p)H_{\mathbb{B}}\left(1 - e^{-|S_0+l|^2}\right) - pH_{\mathbb{B}}\left(1 - e^{-|S_1+l|^2}\right), \end{aligned} \quad (4.33)$$

where l is the feedback control signal of the Dolinar receiver. Note that here we assume the binary output of the photon counter, i.e., it only distinguishes whether or not any positive number of photons arrives. This assumption may not hurt the resulting information rate, under the assumption of average photon number per symbol, $\mathcal{E} \rightarrow 0$. We assume that $0 < p < 1/2$. Since we are highly limited in terms of the average photon number \mathcal{E} while having freedom to choose the feedback control l without any bound, the mean of the optimum input amplitudes $\{S_0, S_1\}$ should be 0. Therefore,

$\{S_0, S_1\}$ should satisfy the following two equations,

$$\begin{aligned} (1-p) \cdot S_0 + p \cdot S_1 &= 0, \\ (1-p) \cdot |S_0|^2 + p \cdot |S_1|^2 &= \mathcal{E}, \end{aligned} \quad (4.34)$$

whose solution can be written as

$$S_0^* = -\sqrt{\frac{\mathcal{E} \cdot p}{1-p}}, \quad S_1^* = \sqrt{\frac{\mathcal{E} \cdot (1-p)}{p}}. \quad (4.35)$$

From Lemma 4.1, it was shown that the optimum l that maximizes binary input/output mutual information of the Dolinar receiver is

$$l^* = \frac{(1-p)S_0^* - pS_1^*}{p - (1-p)} = \frac{-2\sqrt{\mathcal{E} \cdot p \cdot (1-p)}}{2p-1}. \quad (4.36)$$

Therefore, the resulting *effective* photon arrival rates of the two hypotheses become

$$\begin{aligned} \lambda_0 &:= |S_0 + l|^2 = \frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2}, \\ \lambda_1 &:= |S_1 + l|^2 = \frac{(1-p) \cdot \mathcal{E}}{p(1-2p)^2}, \end{aligned} \quad (4.37)$$

respectively. If the prior probability of H_1 , p , is a constant number in the range of $(0, 1/2)$, then both $\lambda_0, \lambda_1 \rightarrow 0$, as $\mathcal{E} \rightarrow 0$. However, we will shortly see that the optimum p that maximizes the mutual information depends on \mathcal{E} , and goes to 0. Then the question is how fast p goes to 0, and whether $\lambda_0, \lambda_1 \rightarrow 0$ at the optimum p . To answer this question, we will first write out the mutual information of the binary input/output channel with the Dolinar receiver as a function of the parameter p , the prior probability of H_1 . We will denote the corresponding mutual information as $I(p)$.

$$I(p) = H_{\mathbb{B}} \left(1 - (1-p)e^{-\lambda_0} + pe^{-\lambda_1} \right) - (1-p)H_{\mathbb{B}}(1 - e^{-\lambda_0}) - pH_{\mathbb{B}}(1 - e^{-\lambda_1}) \quad (4.38)$$

for λ_0 and λ_1 in (4.37).

Now let us assume that $\lambda_0, \lambda_1 \rightarrow 0$ as $\mathcal{E} \rightarrow 0$. Under this assumption, we will find the optimum p^* that maximizes $I(p)$, and then will verify this assumption. When $\lambda_0, \lambda_1 \rightarrow 0$,

$$\begin{aligned}
1 - (1-p)e^{-\lambda_0} - pe^{-\lambda_1} &= ((1-p)\lambda_0 + p\lambda_1) \left(1 - \frac{1}{2} \cdot \frac{(1-p)\lambda_0^2 + p\lambda_1^2}{(1-p)\lambda_0 + p\lambda_1} \right) + O(\lambda_0^3 + p\lambda_1^3) \\
&= \frac{\mathcal{E}}{(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{(p^3 + (1-p)^3)\mathcal{E}}{p(1-p)(1-2p)^2} \right) + O\left(\frac{\mathcal{E}^3}{p^2}\right), \\
1 - e^{-\lambda_0} &= \lambda_0 - \frac{1}{2}\lambda_0^2 + O(\lambda_0^3) \\
&= \frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2} \right) + O(p^3\mathcal{E}^3), \\
1 - e^{-\lambda_1} &= \lambda_1 - \frac{1}{2}\lambda_1^2 + O(\lambda_1^3) \\
&= \frac{(1-p) \cdot \mathcal{E}}{p(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{(1-p)\mathcal{E}}{p(1-2p)^2} \right) + O\left(\frac{\mathcal{E}^3}{p^3}\right).
\end{aligned} \tag{4.39}$$

By using these approximations and additionally, $H_B(x) = -x \log x + x + O(x^2)$ as $x \rightarrow 0$,

$$\begin{aligned}
&H_B\left(1 - (1-p)e^{-\lambda_0} - pe^{-\lambda_1}\right) \\
&= -\frac{\mathcal{E}}{(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{(p^3 + (1-p)^3)\mathcal{E}}{p(1-p)(1-2p)^2} \right) \\
&\quad \times \log \left(\frac{\mathcal{E}}{(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{(p^3 + (1-p)^3)\mathcal{E}}{p(1-p)(1-2p)^2} \right) \right) \\
&\quad + \frac{\mathcal{E}}{(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{(p^3 + (1-p)^3)\mathcal{E}}{p(1-p)(1-2p)^2} \right) + O\left(\frac{\mathcal{E}^3}{p^2} \log \mathcal{E} + \mathcal{E}^2\right), \\
&H_B\left(1 - e^{-\lambda_0}\right) \\
&= -\frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2} \right) \\
&\quad \times \log \left(\frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2} \right) \right) \\
&\quad + \frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{p \cdot \mathcal{E}}{(1-p)(1-2p)^2} \right) + O\left(p^3 \cdot \mathcal{E}^3 \log(p \cdot \mathcal{E}) + (p \cdot \mathcal{E})^2\right),
\end{aligned}$$

$$\begin{aligned}
& H_{\text{B}}(1 - e^{-\lambda_1}) \\
&= -\frac{(1-p) \cdot \mathcal{E}}{p(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{(1-p)\mathcal{E}}{p(1-2p)^2}\right) \log \left(\frac{(1-p) \cdot \mathcal{E}}{p(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{(1-p)\mathcal{E}}{p(1-2p)^2}\right)\right) \\
&\quad + \frac{(1-p) \cdot \mathcal{E}}{p(1-2p)^2} \left(1 - \frac{1}{2} \cdot \frac{(1-p)\mathcal{E}}{p(1-2p)^2}\right) + O\left(\frac{\mathcal{E}^3}{p^3} \log \frac{\mathcal{E}}{p} + \frac{\mathcal{E}^2}{p^2}\right).
\end{aligned}$$

Therefore,

$$\begin{aligned}
I(p) &= H_{\text{B}}\left(1 - (1-p)e^{-\lambda_0} + pe^{-\lambda_1}\right) - (1-p)H_{\text{B}}(1 - e^{-\lambda_0}) - pH_{\text{B}}(1 - e^{-\lambda_1}) \\
&= \left(\frac{\mathcal{E}}{1-2p} - \frac{\mathcal{E}^2}{2} \cdot \frac{(1-p+p^2)}{p(1-p)(1-2p)^3}\right) \log \frac{1-p}{p} + O\left(\frac{\mathcal{E}^3}{p^2} \log \frac{\mathcal{E}}{p} + \frac{\mathcal{E}^2}{p}\right) \\
&= \left(\frac{\mathcal{E}}{1-2p} - \frac{\mathcal{E}^2}{2} \cdot \frac{1+6p}{p}\right) \log \frac{1}{p} + O\left(\frac{\mathcal{E}^3}{p^2} \log \frac{\mathcal{E}}{p} + \frac{\mathcal{E}^2}{p}\right).
\end{aligned} \tag{4.40}$$

The derivative of $I(p)$ can then be approximated as

$$\frac{d}{dp} I(p) \approx -\frac{\mathcal{E}}{p} + \frac{\mathcal{E}^2}{2p^2} \log \frac{1}{p}, \tag{4.41}$$

and it can be checked that when $p = \frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}$, it satisfies $\frac{d}{dp} I(p) \approx \frac{-2 \log \log \frac{1}{\mathcal{E}}}{(\log \mathcal{E})^2} \rightarrow 0$ as $\mathcal{E} \rightarrow 0$. At $p = \frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}$, we can also validate that $\lambda_0, \lambda_1 \rightarrow 0$ as $\mathcal{E} \rightarrow 0$, i.e.,

$$\begin{aligned}
\lambda_0 &= \frac{\mathcal{E}^2 \log \frac{1}{\mathcal{E}}}{2 \left(1 - \frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}\right) \left(1 - \frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}\right)^2} \rightarrow 0, \\
\lambda_1 &= \frac{2 \left(1 - \frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}\right)}{\log \frac{1}{\mathcal{E}} \left(1 - \frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}\right)^2} \rightarrow 0.
\end{aligned} \tag{4.42}$$

When we plug the optimum $p^* = \frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}$ into (4.40),

$$\begin{aligned}
I(p^*) &= \left(\frac{\mathcal{E}}{1 - \mathcal{E} \log \frac{1}{\mathcal{E}}} - \mathcal{E} \cdot \frac{1 + 3\mathcal{E} \log \frac{1}{\mathcal{E}}}{\log \frac{1}{\mathcal{E}}}\right) \left(\log \frac{2}{\mathcal{E}} - \log \log \frac{1}{\mathcal{E}}\right) + O\left(\frac{\mathcal{E}}{\log \frac{1}{\mathcal{E}}}\right) \\
&= \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}).
\end{aligned} \tag{4.43}$$

Therefore, (4.43) is the maximum mutual information of the binary coherent states

$\{|S_0\rangle, |S_1\rangle\}$ under the average photon number constraint of \mathcal{E} , when the received coherent state is measured by the coherent receiver.

Now, we generalize this result to M -ary input states. We will use mathematical induction to show that we cannot exceed the mutual information in (4.43) even when we generalize the inputs to M -ary coherent states, satisfying the same average photon number constraint of \mathcal{E} . Let us define $R_M(\mathcal{E})$ as the maximum mutual information of the optical channel with the coherent receiver that detects a state drawn from the M -ary input states under the average photon number of \mathcal{E} . In (4.43), we showed that

$$R_2(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}). \quad (4.44)$$

Now, we assume that for a finite M ,

$$R_M(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}), \quad (4.45)$$

and then we will prove that

$$R_{M+1}(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}). \quad (4.46)$$

By mathematical induction, this proof will imply that for any *finite* M , (4.45) is true, and thus the achievable photon efficiency with coherent receivers is not significantly different from that of direct detection receivers.

First, for a given feedback control signal l , without loss of generality, we arrange $(M + 1)$ -input states $\{|S_1\rangle, |S_2\rangle, \dots, |S_{M+1}\rangle\}$, such that it satisfies

$$|S_0 + l|^2 \leq |S_1 + l|^2 \leq \dots \leq |S_{M+1} + l|^2. \quad (4.47)$$

We denote the prior probabilities of the input states as $\{p_1, p_2, \dots, p_{M+1}\}$. The binary

output distribution of the channel becomes

$$p_Y = \left\{ \sum_{i=1}^{M+1} p_i \cdot e^{-|S_i+l|^2}, 1 - \left(\sum_{i=1}^{M+1} p_i \cdot e^{-|S_i+l|^2} \right) \right\}, \quad (4.48)$$

and the resulting mutual information of the M -ary input channel becomes

$$I(H, Y) = H_B \left(\sum_{i=1}^{M+1} p_i \cdot e^{-|S_i+l|^2} \right) - \sum_{i=1}^{M+1} p_i H_B \left(e^{-|S_i+l|^2} \right). \quad (4.49)$$

Now let us define two random variables N_1, N_2 from H as follows.

$$N_1 = \begin{cases} 0, & \text{if } H = 1, \dots, M \\ 1, & \text{o.w.} \end{cases}, \quad N_2 = \begin{cases} H, & \text{if } H = 1, \dots, M \\ 0, & \text{o.w.} \end{cases}.$$

Since H and (N_1, N_2) are bijective, by chain rule,

$$I(H; Y) = I(N_1, N_2; Y) = I(N_1; Y) + I(N_2; Y|N_1). \quad (4.50)$$

Note that

$$\begin{aligned} I(N_2; Y|N_1) &= \left(\sum_{i=1}^M p_i \right) \cdot I(N_2; Y|N_1 = 1) + p_{M+1} \cdot I(N_2; Y|N_1 = 0) \\ &= \left(\sum_{i=1}^M p_i \right) \cdot \left(H_B \left(\sum_{j=1}^M \frac{p_j}{\left(\sum_{i=1}^M p_i \right)} \cdot e^{-|S_j+l|^2} \right) - \sum_{j=1}^M \frac{p_j}{\left(\sum_{i=1}^M p_i \right)} H_B \left(e^{-|S_j+l|^2} \right) \right) \end{aligned} \quad (4.51)$$

since $I(N_2; Y|N_1 = 0) = 0$. The average number of *effective* photons used to encode the information in N_2 , which will be denoted as \mathcal{E}_2 , equals

$$\mathcal{E}_2 = \sum_{j=1}^M \frac{p_j}{\left(\sum_{i=1}^M p_i \right)} |S_j - \bar{S}|^2 \quad (4.52)$$

where $\bar{S} = \frac{\sum_{i=1}^M p_i \cdot S_i}{\sum_{i=1}^M p_i}$, the mean of the M -signal amplitudes. When we calculate the *effective* photon numbers, we subtract \bar{S} from S_j , since we can make any common offset to the signals $\{S_j\}$ by adding the desired amount to l without any cost. From (4.51) and the definition of $R_M(\mathcal{E})$,

$$I(N_2; Y|N_1) \leq \left(\sum_{i=1}^M p_i \right) \cdot R_M(\mathcal{E}_2). \quad (4.53)$$

On the other hand, the mutual information between N_1 and Y becomes

$$\begin{aligned} I(N_1; Y) = & H_B \left(\sum_{i=1}^{M+1} p_i \cdot e^{-|S_i+l|^2} \right) \\ & - \left(\sum_{i=1}^M p_i \right) \cdot H_B \left(\sum_{j=1}^M \frac{p_j}{\left(\sum_{i=1}^M p_i \right)} e^{-|S_j+l|^2} \right) - p_{M+1} \cdot H_B \left(e^{-|S_{M+1}+l|^2} \right). \end{aligned} \quad (4.54)$$

The channel distribution of $p_{Y|N_1}$ is

$$p_{Y|N_1=0} = \begin{cases} \sum_{j=1}^M \frac{p_j}{\left(\sum_{i=1}^M p_i \right)} e^{-|S_j+l|^2} & \text{if } Y = 0 \\ 1 - \sum_{j=1}^M \frac{p_j}{\left(\sum_{i=1}^M p_i \right)} e^{-|S_j+l|^2} & \text{if } Y = 1 \end{cases} \quad (4.55)$$

$$p_{Y|N_1=1} = \begin{cases} e^{-|S_{M+1}+l|^2} & \text{if } Y = 0 \\ 1 - e^{-|S_{M+1}+l|^2} & \text{if } Y = 1 \end{cases} \quad (4.56)$$

where the input distribution $p_{N_1} = \left\{ \sum_{i=1}^M p_i, p_{M+1} \right\}$.

Now let us define a new channel $q_{Y|N_1}$ such that

$$q_{Y|N_1=0} = \begin{cases} e^{-|\bar{S}+l|^2} & \text{if } Y = 0 \\ 1 - e^{-|\bar{S}+l|^2} & \text{if } Y = 1 \end{cases} \quad (4.57)$$

$$q_{Y|N_1=1} = p_{Y|N_1=1}, \quad Y \in \{0, 1\} \quad (4.58)$$

for $\bar{S} = \frac{\sum_{i=1}^M p_i \cdot S_i}{\sum_{i=1}^M p_i}$. The average number of photons for the channel distributions $q_{Y|N_1}$ under the input distribution p_{N_1} is denoted as \mathcal{E}_1 and calculated below:

$$\mathcal{E}_1 = \left(\sum_{i=1}^M p_i \right) \cdot |\bar{S}|^2 + p_{M+1} \cdot |S_{M+1}|^2 \quad (4.59)$$

From the definition of $R_2(\mathcal{E})$,

$$I(q_{Y|N_1}, p_{N_1}) \leq R_2(\mathcal{E}_1). \quad (4.60)$$

Now we will show that

$$I(p_{Y|N_1}, p_{N_1}) \leq I(q_{Y|N_1}, p_{N_1}), \quad (4.61)$$

which will imply $I(N_1; Y)$ in (4.50) is less than

$$I(N_1; Y) \leq R_2(\mathcal{E}_1). \quad (4.62)$$

We will use the following lemma for this step.

Lemma 4.5. *For a binary input/output channel $W_{Y|X}$ with the input distribution $p_X = \{p_0, p_1\}$, let the channel distribution conditioned on $Y = 1$ be $W_{Y|X=1} = \{t_1, 1 - t_1\}$ and on $Y = 0$ be $W_{Y|X=0} = \{t_0, 1 - t_0\}$ while $t_0 \geq t_1$. Then, the mutual information of this channel, parameterized by t_0 , $f(t_0) := I(p_X; W_{Y|X})$, monotonically decreases as t_0 goes to t_1 .*

Proof. Let us denote the channel distribution $W_{Y|X}$ with the parameter t_0 as a matrix $W_{t_0} := \begin{pmatrix} t_0 & 1 - t_0 \\ t_1 & 1 - t_1 \end{pmatrix}$. For $t_0 \geq t'_0 \geq t_1$, $\exists r \in [0, 1)$ s.t. $r \cdot W_{t_0} + (1 - r) \cdot W_{t_1} = W_{t'_0}$. Since mutual information $I(p_X, W_{Y|X})$ is convex in $W_{Y|X}$, $f(t_0)$ is convex in t_0 . Therefore, $r \cdot f(t_0) + (1 - r) \cdot f(t_1) \geq f(t'_0)$. Since $f(t_1) = 0$, the convexity gives $f(t_0) \geq r \cdot f(t_0) \geq f(t'_0)$, i.e., $f(t_0)$ is a monotonically decreasing function in t_0 as $t_0 \geq t_1$ goes closer to t_1 . \square

If the following inequalities hold

$$e^{-|S_{M+1}+l|^2} \leq \sum_{j=1}^M \frac{p_j}{\left(\sum_{i=1}^M p_i\right)} e^{-|S_j+l|^2} \leq e^{-|\bar{S}+l|^2}, \quad (4.63)$$

then by using Lemma 4.5 and the definition of the previously introduced channel distributions $p_{Y|N_1}$ and $q_{Y|N_1}$, we can prove (4.61). In (4.63), the first inequality is valid from our assumption in (4.47). The second inequality is also valid since $e^{-|x+l|^2}$ is concave in x when $|x+l|^2 \leq 1/2$, and we are interested in the regime of $|x+l|^2 \rightarrow 0$. Therefore, we have shown that (4.62) is true. By combining (4.62) and (4.53),

$$I(H; Y) = I(N_1; Y) + I(N_2; Y|N_1) \leq R_2(\mathcal{E}_1) + \left(\sum_{i=1}^M p_i\right) \cdot R_M(\mathcal{E}_2) \quad (4.64)$$

where \mathcal{E}_1 and \mathcal{E}_2 are (4.59) and (4.52), respectively. Also note that

$$\begin{aligned} \mathcal{E}_1 + \left(\sum_{i=1}^M p_i\right) \cdot \mathcal{E}_2 &= \left(\sum_{i=1}^M p_i\right) \cdot |\bar{S}|^2 + p_{M+1} \cdot |S_{M+1}|^2 + \sum_{i=1}^M p_i \cdot |S_i - \bar{S}|^2 \\ &= \left(\sum_{i=1}^M p_i\right) \cdot (2 \cdot |\bar{S}|^2 + |S_i|^2 - 2S_i\bar{S}) + p_{M+1} \cdot |S_{M+1}|^2 \\ &= \sum_{i=1}^{M+1} p_i |S_i|^2 = \mathcal{E}. \end{aligned} \quad (4.65)$$

The second equality holds since $2 \left(\sum_{i=1}^M p_i S_i\right) \bar{S} = 2 \left(\sum_{i=1}^M p_i\right) \bar{S}^2$ from the definition of \bar{S} .

Therefore, when we denote $\mathcal{E}_1 = (1 - \alpha) \cdot \mathcal{E}$ and $\mathcal{E}_2 = \alpha \cdot \mathcal{E} / \beta$ for $\beta := \left(\sum_{i=1}^M p_i\right) < 1$ and some $\alpha \in (0, 1)$, the upper bound of $I(H; Y)$ in (4.64) becomes

$$I(H; Y) \leq R_2((1 - \alpha) \cdot \mathcal{E}) + \beta \cdot R_M(\alpha \cdot \mathcal{E} / \beta). \quad (4.66)$$

From (4.44) and the assumption (4.45),

$$\begin{aligned}
I(H; Y) &\leq ((1 - \alpha)\mathcal{E}) \log \frac{1}{((1 - \alpha)\mathcal{E})} - ((1 - \alpha)\mathcal{E}) \log \log \frac{1}{((1 - \alpha)\mathcal{E})} \\
&\quad + \beta \cdot \left((\alpha \cdot \mathcal{E}/\beta) \log \frac{1}{(\alpha \cdot \mathcal{E}/\beta)} - (\alpha \cdot \mathcal{E}/\beta) \log \log \frac{1}{(\alpha \cdot \mathcal{E}/\beta)} \right) + O(\mathcal{E}) \quad (4.67) \\
&= \log \frac{1}{\mathcal{E}} - \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}).
\end{aligned}$$

This inequality holds for every $(M+1)$ -ary input states with $p_i > 0$, for $i = 1, \dots, M+1$, under the average photon number constraint of \mathcal{E} . Therefore,

$$R_{M+1}(\mathcal{E}) \leq \log \frac{1}{\mathcal{E}} - \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}), \quad (4.68)$$

and by induction, (4.45) holds for any finite M . This concludes the proof of Theorem 4.4.

Superadditivity of Quantum Channel Coding Rate with Finite Blocklength Quantum Measurements

■ 5.1 Background

How many classical bits per channel use can be reliably communicated over a quantum channel? This has been a central question in quantum information theory in an effort to understand the intrinsic limit on the classical capacity of physical quantum channels such as the optical fiber or free-space optical channels. The Holevo limit of a quantum channel is an upper bound to the Shannon capacity of the classical channel induced by pairing the quantum channel with any specific receiver measurement [20, 23]. The Holevo limit is in principle also an achievable information rate, and is known for several important practical channels, such as the lossy bosonic channel [14]. However, a receiver that attains the Holevo capacity, must in general make joint (*collective*) measurements over long codeword blocks. Such measurements cannot be realized by detecting single modulation symbols followed by classical post processing. We call this phenomenon of a joint-detection receiver (JDR) being able to yield a higher information rate (in

error-free bits communicated per channel use) than what is possible by any single-symbol receiver measurement, as *superadditivity* of capacity¹. There are several JDR measurements that are known to achieve the Holevo capacity—the square-root measurement (SRM) [20], Helstrom’s minimum probability of error measurement [21], the sequential-decoding measurement [15, 47], the successive-cancellation decoder for the quantum polar code [17, 46], and a two-stage near-unambiguous-detection receiver [42]. There are a few characteristics that are common to each one of these measurements. First, the size of the joint-detection measurement is tied to the blocklength of the code, i.e., the measurement must act on the entire codeword and hence its size must increase with the length of the codeword. Second, none of these measurement specifications translate readily to a realizable receiver in the context of optical communication. Since it is known that a simple laser-light (coherent-state) modulation achieves the Holevo capacity of the lossy bosonic channel [14], almost all the complexity in achieving the ultimate limit to the reliable communication rate lies at the receiver. Finally, none of these capacity-achieving measurements tell us how the achievable information rate increases with the size of the receiver measurement. Since the complexity of implementing a joint quantum measurement over N channel symbols in general grows exponentially with N , it is of great practical interest to find how the maximum achievable information rate (error-free bits per channel use) scales with the size N of the joint-detection receiver (while imposing no constraint whatsoever on the classical code complexity). In this chapter, we shed some light on this, for classical communication over a quantum channel using a pure-state alphabet—the so-called pure-state classical-quantum (cq) channel (the lossy bosonic channel being an example)—by proving a general lower

¹We would like to clarify that the more prevalent use of the term *superadditivity* of capacity refers to the scenario when a quantum channel has a higher classical communication capacity when using transmitted states that are entangled over multiple channel uses [19]. For the bosonic channel, it was shown that entangled inputs at the transmitter cannot get a higher capacity [14]. However, one *can* get a higher capacity on the bosonic channel—as compared to what is possible by any optical receiver that measures one channel output at a time—by using *entangling* (or joint-detection) measurements at the receiver. As the number of symbols over which the receiver acts increases, the capacity steadily increases. In this dissertation, we use the term *superadditivity* in this latter context, and provide a general bound on the scaling of the capacity with the length of the joint measurement. This usage of the term was first adopted by Sasaki, et. al. [36].

bound on the finite-measurement-length capacity.

Finally, we would like to remark on an important difference between our results in this chapter and the finite-blocklength rate over a cq-channel derived using second-order asymptotics (channel dispersion) [30, 43]. The latter techniques explore how the achievable rate R_N/N (bits per channel use), at a *given decoding error threshold* ϵ , increases when both the code length and the measurement length increase *together*. We consider the *asymptotic* capacity C_N/N (error-free bits per channel use), while imposing a constraint on the receiver to make collective measurements over N channel outputs, but with no restriction on the complexity of the classical outer code.

■ 5.2 Introduction and Problem Statement

The classical capacity of a quantum channel is defined as the maximum number of information bits that can be modulated into the input quantum states and reliably decoded at the receiver with a set of quantum measurements as the number of transmissions N_c goes to infinity. Consider a pure-state classical-quantum (cq) channel $W : x \rightarrow |\psi_x\rangle$, where $x \in \mathcal{X}$ is the classical input, and $\{|\psi_x\rangle\} \in \mathcal{H}$ are corresponding modulation symbols at the output of the channel. One practical example of a pure-state cq channel is the single-mode lossy optical channel $\mathcal{N}_\eta : \alpha \rightarrow |\sqrt{\eta}\alpha\rangle$, where $\alpha \in \mathbb{C}$ is the complex field amplitude at the input of the channel, $\eta \in (0, 1]$ is the transmissivity (the fraction of input power that appears at the output), and $|\sqrt{\eta}\alpha\rangle$ is the quantum description of an ideal laser-light pulse, a coherent-state².

In refs. [20, 23], it was shown that the classical capacity of a cq channel W is given by

$$C = \max_{P_X} \text{Tr}(-\rho \log \rho), \quad (5.1)$$

²It is important to note here the difference between a classical channel and a classical-quantum channel. There is no physical measurement that can noiselessly measure the output amplitude $\sqrt{\eta}\alpha$. Any specific choice of an optical receiver—such as homodyne, heterodyne or direct-detection—induces a specific discrete memoryless *classical* channel $p(\beta|\alpha)$ between the input α and the measurement result β . The Shannon capacity of this classical channel, for any given measurement choice, is strictly smaller than the Holevo capacity of the cq channel \mathcal{N}_η .

where $\rho = \sum_{x \in \mathcal{X}} P_X(x) |\psi_x\rangle\langle\psi_x|$. The states $|\psi_x\rangle$, $x \in \mathcal{X}$, are normalized vectors in a complex Hilbert space \mathcal{H} , $\langle\psi_x|$ is the Hermitian conjugate vector of $|\psi_x\rangle$, and ρ is a *density operator*, a linear combination of the outer products $|\psi_x\rangle\langle\psi_x|$ with weights $P_X(x)$. The capacity can also be written as $C = \max_{P_X} S(\rho)$, where $S(\rho) = \text{Tr}(-\rho \log \rho)$ is the von Neumann entropy of the density operator ρ .

For an input codeword (x_1, \dots, x_{N_c}) , the received codeword is a tensor product state, $|\psi_{x_1}\rangle \otimes \dots \otimes |\psi_{x_{N_c}}\rangle$, which is jointly detected by an orthogonal projective measurement in the N_c -fold Hilbert space $\mathcal{H}^{\otimes N_c}$. When the received codeword is projected into the orthogonal measurement vectors $\{|\Phi_k\rangle\}$, $k \in \mathcal{K}$, which resolve the identity, i.e., $\sum_k |\Phi_k\rangle\langle\Phi_k| = \mathbb{1}$, in $\mathcal{H}^{\otimes N_c}$, the classical output k is observed with probability equal to the magnitude squared of the inner product between the received codeword state, $|\psi_{x_1}\rangle \otimes \dots \otimes |\psi_{x_{N_c}}\rangle$, and the measurement vector $|\Phi_k\rangle$ corresponding to the output k . The orthogonal projective measurement is designed to decode the received codewords with as small error probability as possible. For any rate $R < C$, a block code of length N_c and rate R , generated by picking each of the N_c symbols of $e^{N_c R}$ codewords randomly from the distribution P_X^* that attains the maximum in Eq. (5.1), paired with an appropriate quantum measurement acting jointly on the received codeword in $\mathcal{H}^{\otimes N_c}$, can attain an arbitrarily small probability of error as $N_c \rightarrow \infty$ [15, 20, 23, 42].

To achieve this capacity, however, a joint detection receiver (JDR) needs to be implemented, which can measure the length- N_c sequence of states jointly and decode it reliably among $e^{N_c C}$ possible messages. The number of measurement outcomes thus scales exponentially with the length of the codeword N_c . Hence, the complexity of physical implementation (in terms of number of elementary finite-length quantum operations) of the receiver in general also grows exponentially with N_c . Considering this exponential growth in complexity, one might want to limit the maximum length $N \leq N_c$ of the sequence of states to be jointly detected at the receiver, independent of the length of the codeword N_c . However, there is no guarantee that such quantum measurements of fixed blocklengths can still achieve the ultimate capacity of the quantum channel.

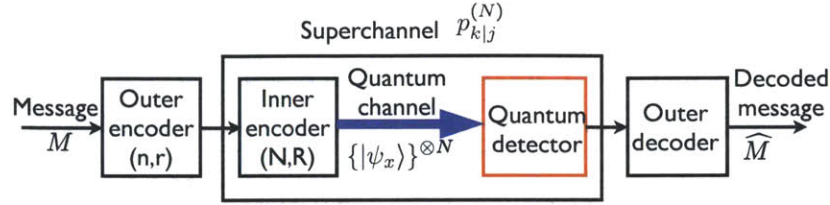


Figure 5.1. Concatenated coding over a classical-quantum channel

Our goal in this chapter is to study the trade-off between information rate and receiver complexity, for classical communication over a quantum channel. We investigate the maximum number of classical information bits that can be reliably decoded per use of the quantum channel, when quantum states of a finite blocklength $N \leq N_c$ are jointly measured, with no restriction on the complexity of the overall classical code ($N_c \rightarrow \infty$). As the number of channel outputs N jointly measured increases, the maximum number of classical bits extracted per use of the channel increases. We call this phenomenon *superadditivity* of the maximum achievable classical information rate over a quantum channel. After the receiver detects the quantum states, it can collect all the classical information extracted from each block of length N , and then apply any classical decoding algorithm over the collected information to decode the transmitted message reliably. To explain how it works, we introduce the architecture of concatenation over a quantum channel in Fig. 5.1.

In the communication system depicted in Fig. 5.1, a *concatenated code* is used to transmit the message M over the quantum channel. For an inner code of length N and rate R , there can be a total of e^{NR} inputs to the inner encoder, $J \in \{1, \dots, e^{NR}\}$. The inner encoder maps each input J to a length- N classical codeword, which maps to a length- N sequence of quantum states at the output of the quantum channel. The quantum joint-detection receiver measures the length- N quantum codeword and generates an estimate $K \in \{1, \dots, e^{NR}\}$ of the encoded message J^3 . For a good inner

³Note that, in general, in order to attain the maximum mutual information of the superchannel induced by the inner code-decoder pair, the number of outputs of the inner decoder may need to be greater than e^{NR} [9, 40]. But we will stick to the case when the output of the inner decoder makes a hard decision on the input message.

code and joint measurement with N large, the estimate would generally match the input message. But for a fixed N , the error probability may not be close to 0. The inner encoder, the classical-quantum channel, and the inner decoder (the joint detection receiver), collectively forms a discrete memoryless *superchannel*, with transition probabilities $p_{k|j}^{(N)} := \Pr(K = k|J = j)$. We define the maximum mutual information of this superchannel, over all choices of inner codes of blocklength N , and over all choices of inner-decoder joint measurements of length N as:

$$C_N := \max_{p_j} \max_{\{N\text{-symbol inner code-measurement pairs}\}} I(p_j, p_{k|j}^{(N)}). \quad (5.2)$$

A classical Shannon-capacity-achieving outer code can be used to reliably communicate information through the superchannel of the maximum mutual information C_N . When an outer code of length n and rate r is adopted, the total number of messages transmitted by the code is $e^{nr} = e^{nN(r/N)}$. Since the overall length of the concatenated code, which is composed of the inner code and the outer code, is $N_c = nN$, the total rate of the concatenated code is $R_c = r/N$. By Shannon's coding theorem, for any rate $r < C_N$, there exists an outer code of length n and rate r such that the decoding error can be made arbitrarily small as $n \rightarrow \infty$. Therefore, the maximum information rate achievable by the concatenated code *per use of the quantum channel* can approach C_N/N .

From the definition of C_N , superadditivity of the quantity, i.e., $C_{N_1} + C_{N_2} \leq C_{N_1+N_2}$, can be shown. This implies the existence of the limit $\lim_{N \rightarrow \infty} C_N/N$. Holevo [24] showed that the limit is equal to the ultimate capacity of the quantum channel, $\lim_{N \rightarrow \infty} C_N/N = C = \max_{P_X} S(\rho)$. Therefore, C_N/N is an increasing sequence in N with its limit equal to the capacity.

The question we want to answer is: *How does the maximum achievable information rate C_N/N increase as the length of the quantum measurement, N , increases?*

On the receiver side, since quantum processing occurs only at the quantum decoder for the inner code of a finite length N , the complexity of the quantum processing only

depends on N , but not on the outer code length n . Therefore, the trade-off between the (rate) performance and the (quantum) complexity can be captured by how fast C_N/N increases with N . It is known that for some examples of input states, strict superadditivity of C_N can be demonstrated [16, 24, 33]. However, the calculation of C_N , even for a pure-state binary alphabet, is extremely hard for $N > 1$ because the complexity of optimization increases exponentially with N . The concatenated coding scheme described above was considered in [18] for the lossy bosonic channel, where some examples of inner codes and structured optical joint-detection receivers were found for which $C_N/N > C_1$ holds for a binary coherent-state modulation alphabet.

Instead of aiming to calculate the exact C_N , in this chapter, a lower bound of C_N/N , which becomes tight for large enough N , will be derived. From this bound, it will be possible to calculate the inner code blocklength N at which a given fraction of the ultimate capacity is achievable. A new framework for understanding the strict superadditivity of C_N in quantum channels will also be provided, which is different from the previous explanation of the phenomenon by *entangled measurements* and the resulting memory in the quantum channel [35].

The rest of the chapter is outlined as follows. In Section 5.3, examples of quantum channels where strict superadditivity $C_1 < C$ holds, will be demonstrated. The main theorem that states a lower bound on C_N/N will be given in Section 5.4 with examples to show how to use the main theorem to calculate a blocklength N to achieve a given fraction of the capacity. The theorem will be proved in Section 5.5. Thereafter in Section 5.6, the effect of superadditivity due to finite-blocklength inner-code measurements in a concatenated coding architecture will be compared between a quantum channel and a classical discrete memoryless channel. An approximation of the lower bound of C_N/N will also be provided by introducing a quantum version of *channel dispersion* V , with a unifying picture encompassing quantum and classical channels. Section 5.7 will conclude this chapter.

■ 5.3 Strict Superadditivity of C_N

Before investigating how C_N/N increases with N , we will show examples where strict superadditivity of C_N can be shown, i.e., $C_1 < C$. As discussed before, given a set of input states $\{|\psi_x\rangle\}$, $x \in \mathcal{X}$, C can be calculated from Holevo's result by finding the optimum input distribution that maximizes the von Neumann entropy $S(\rho) = \text{Tr}(-\rho \log \rho)$ where $\rho = \sum_x P_X(x) |\psi_x\rangle\langle\psi_x|$. Calculating C_1 on the other hand requires finding a set of measurements as well as an input distribution to maximize the resulting mutual information, where the measurement acts on one channel symbol at a time. For general input states, this is a hard optimization problem, since the measurement that maximizes C_1 may not be a projective measurement, and could be a Positive Operator Valued Measure (POVM)—the most general description of a quantum measurement—and furthermore the optimum POVM could have up to $|\mathcal{X}|(|\mathcal{X}| + 1)/2$ outcomes [9, 40]. However, for binary pure states $\{|\psi_0\rangle, |\psi_1\rangle\}$ of inner product $|\langle\psi_0|\psi_1\rangle| = \gamma$, C_1 and C can be calculated as simple functions of the inner product $\gamma = |\langle\psi_0|\psi_1\rangle|$, and strict superadditivity can be shown, as summarized below [24].

The first step to calculate C for the binary input channel is to find the eigenvalues of ρ under an input distribution $\{1 - q, q\}$. For $\rho = (1 - q)|\psi_0\rangle\langle\psi_0| + q|\psi_1\rangle\langle\psi_1|$, the eigenvectors of ρ have a form of $|\psi_0\rangle + \beta|\psi_1\rangle$ with some β that satisfies

$$\rho(|\psi_0\rangle + \beta|\psi_1\rangle) = \sigma(|\psi_0\rangle + \beta|\psi_1\rangle) \quad (5.3)$$

with eigenvalues σ . By solving the equation, we obtain the two eigenvalues as:

$$\begin{aligned} \sigma_1 &= \frac{1}{2} \left(1 - \sqrt{1 - 4q(1 - q)(1 - \gamma^2)} \right), \text{ and} \\ \sigma_2 &= \frac{1}{2} \left(1 + \sqrt{1 - 4q(1 - q)(1 - \gamma^2)} \right), \end{aligned} \quad (5.4)$$

and the resulting von Neumann entropy,

$$S(\rho) = \text{Tr}(-\rho \log \rho) = -\sigma_1 \log \sigma_1 - \sigma_2 \log \sigma_2, \quad (5.5)$$

where $|\langle\psi_0|\psi_1\rangle| = \gamma$. From this equation, it can be shown that $S(\rho)$ for the binary inputs is maximized at $q = 1/2$, and the resulting capacity of the binary channel,

$$C = \max_{P_X} S(\rho) = -\frac{1-\gamma}{2} \log \frac{1-\gamma}{2} - \frac{1+\gamma}{2} \log \frac{1+\gamma}{2}. \quad (5.6)$$

For the binary channel, C_1 is attained by the equiprior input distribution and a binary-valued projective measurement in the span of $\{|\psi_0\rangle, |\psi_1\rangle\}$ —the same measurement that minimizes the average error probability of discriminating between equally-likely states $|\psi_0\rangle$ and $|\psi_1\rangle$. The derivation of C_1 for the binary case can be found in [24], and is given by:

$$C_1 = \frac{1 - \sqrt{1 - \gamma^2}}{2} \log \left(1 - \sqrt{1 - \gamma^2}\right) + \frac{1 + \sqrt{1 - \gamma^2}}{2} \log \left(1 + \sqrt{1 - \gamma^2}\right). \quad (5.7)$$

The capacity C is strictly greater than C_1 for all $0 < \gamma < 1$, which demonstrates the strict superadditivity of C_N for all binary input quantum channels.

In the rest of this section, we will consider the superadditivity of C_N in quantum channels *with an input constraint*, in the context of optical communication. The constraint will be the mean energy of input states. A *coherent state* $|\alpha\rangle$ is the quantum description of a single spatio-temporal-polarization mode of a classical optical-frequency electromagnetic (ideal laser-light) field, where $\alpha \in \mathbb{C}$ is the complex amplitude, and $|\alpha|^2$ is the mean photon number of the mode. Since the energy of a *photon* with angular frequency ω is $E = \hbar\omega$ with $\hbar = h/2\pi$ where the Planck constant $h = 6.63 \times 10^{-34} \text{m}^2\text{kg/s}$, the average energy (in Joules) of the coherent state $|\alpha\rangle$ of a quasi-monochromatic field mode of center frequency ω , is $\hbar|\alpha|^2\omega$. Note that the mean photon number $|\alpha|^2$ is a dimensionless quantity. Therefore, for quasi-monochromatic propagation at a fixed center frequency ω , an average energy constraint on the input states (or equivalently, an average power constraint with a fixed time-slot width) can be represented as a constraint on the average photon number per transmitted state. For example, if the modulation constellation comprises of the set of input states $\{|\alpha_1\rangle, |\alpha_2\rangle, \dots, |\alpha_K\rangle\}$, an average en-

ergy constraint $\hbar\omega\mathcal{E}$ per symbol transmission can be expressed as a constraint on the prior distribution $\{p_1, \dots, p_K\}$, with

$$\sum_{i=1}^K p_i |\alpha_i|^2 \leq \mathcal{E}, \quad (5.8)$$

where \mathcal{E} is the maximum mean photon number per symbol.

The important question of how many bits can be reliably communicated per use (i.e., per transmitted mode) of a pure-loss optical channel of power transmissivity $\eta \in (0, 1]$, under the constraint on the average photon number per transmitted mode \mathcal{E} , was answered in [14]. It was also shown that product coherent state inputs are sufficient to achieve the Holevo capacity of this quantum channel. Since a coherent state $|\alpha\rangle$ of mean photon number $\mathcal{E} = |\alpha|^2$ transforms into another coherent state $|\sqrt{\eta}\alpha\rangle$ of mean photon number $\eta\mathcal{E}$ over the lossy channel, we will henceforth, without loss of generality, subsume the channel loss in the energy constraint, and pretend that we have a lossless channel ($\eta = 1$) with a mean photon number constraint $\mathbb{E}[|\alpha|^2] \leq \mathcal{E}$ per mode (or per ‘channel use’). The capacity of this channel is given by [14]

$$C(\mathcal{E}) = (1 + \mathcal{E}) \log(1 + \mathcal{E}) - \mathcal{E} \log \mathcal{E} \text{ [nats/mode]}, \quad (5.9)$$

and it is achievable with a coherent-state random code with the amplitude α chosen from a circulo-complex Gaussian distribution with variance \mathcal{E} , $p(\alpha) = \exp[-|\alpha|^2/\mathcal{E}]/(\pi\mathcal{E})$.

The number of information bits that can be reliably communicated *per photon*—the photon information efficiency (PIE)—under a mean photon number constraint per mode, \mathcal{E} , is given by $C(\mathcal{E})/\mathcal{E}$ (nats/photon). From (5.9), it can be shown that in order to achieve high PIE, \mathcal{E} must be small. In the $\mathcal{E} \rightarrow 0$ regime, the capacity (5.9) can be approximated as

$$C(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} + \mathcal{E} + o(\mathcal{E}), \quad (5.10)$$

which shows that $\text{PIE} \sim -\log \mathcal{E}$ for $\mathcal{E} \ll 1$. Thus there is no upper limit in principle to

the photon information efficiency.

We will now show that in the high-PIE (low photon number) regime, this ultimate capacity is achievable closely even with a simple Binary Phase Shift Keying (BPSK) coherent state constellation $\{|\sqrt{\mathcal{E}}\rangle, |-\sqrt{\mathcal{E}}\rangle\}$, which satisfies the energy constraint with any prior distribution. The inner product between the two coherent states, $\gamma = |\langle\boldsymbol{\alpha}|\boldsymbol{\beta}\rangle| = \exp[-|\boldsymbol{\alpha} - \boldsymbol{\beta}|^2/2]$. Therefore,

$$|\langle\sqrt{\mathcal{E}}|-\sqrt{\mathcal{E}}\rangle| = \exp[-2\mathcal{E}]. \quad (5.11)$$

By plugging $\gamma = \exp[-2\mathcal{E}]$ into (5.6), we obtain the capacity of the BPSK input constellation,

$$C_{\text{BPSK}}(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} + \mathcal{E} + o(\mathcal{E}), \quad (5.12)$$

which is equal to $C(\mathcal{E})$ for the first- and second-order terms in the limit $\mathcal{E} \rightarrow 0$.

We now ask, for binary coherent-state inputs under the same constraint on the mean photon number per mode \mathcal{E} , how high an information rate is achievable when each mode is detected one-by-one, i.e., $N = 1$. The maximum capacity of the bosonic channel with a binary-input with mean photon number constraint \mathcal{E} , and a $N = 1$ measurement, will be denoted as $C_{1,\text{Binary}}(\mathcal{E})$. For BPSK input states, by using (5.7), the maximum achievable rate at $N = 1$ is

$$C_{1,\text{BPSK}}(\mathcal{E}) = 2\mathcal{E} + o(\mathcal{E}) \quad (5.13)$$

Thus, PIE of the BPSK channel caps off at 2 nats/photon for $N = 1$, while for N large, achievable PIE $\rightarrow \infty$ as $\mathcal{E} \rightarrow 0$.

$C_{1,\text{Binary}}(\mathcal{E})$ can be calculated in the regime $\mathcal{E} \rightarrow 0$ by finding the optimum binary

inputs $\{|\alpha_0\rangle, |\alpha_1\rangle\}$ with distribution $\{1 - q, q\}$ that satisfies the average constraint,

$$(1 - q)|\alpha_0|^2 + q|\alpha_1|^2 \leq \mathcal{E}. \quad (5.14)$$

The following lemma summarizes the result.

Lemma 5.1. *The optimum binary inputs for $N = 1$, are $\alpha_0 = \sqrt{\mathcal{E} \cdot q^*/(1 - q^*)}$ and $\alpha_1 = -\sqrt{\mathcal{E} \cdot (1 - q^*)/q^*}$ with*

$$q^* = \frac{\mathcal{E}}{2} \log \frac{1}{\mathcal{E}}, \quad (5.15)$$

and the resulting $C_{1,\text{Binary}}(\mathcal{E})$ is

$$C_{1,\text{Binary}}(\mathcal{E}) = \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}). \quad (5.16)$$

Proof. Appendix 5.A. □

Compared to the ultimate capacity (5.10) with the same energy constraint, the first-order term of $C_{1,\text{Binary}}(\mathcal{E})$ is the same as that of $C(\mathcal{E})$. But, the difference in the second-order term shows how much less capacity is achievable at $N = 1$ even with the optimized input states. In [6], we showed that (5.16) can be achieved using an on-off keying modulation $\{|\mathbf{0}\rangle, |\alpha\rangle\}$ and a simple on-off direct-detection (photon counting) receiver. Therefore, in the context of optical communication in the high-PIE regime, all of the performance gain from the complex quantum processing in the JDR is captured by the difference between the second-order terms of Eqs. (5.10) and (5.16). In the low photon number regime, this difference in the second-order term can have a significant impact on the practical design of an optical communication system [18]. It would therefore be interesting to ask how large a JDR length N is needed to bridge the gap in the second-order term. To answer this question, we will need the general lower bound on C_N that we develop in the following section.

■ 5.4 Lower Bound on C_N

In this section, a lower bound is derived for the maximum achievable information rate at a finite blocklength N of quantum measurements. By using the result, it will be possible to calculate a blocklength N at which a given fraction of the Holevo capacity of a pure-state cq channel can be achieved. Therefore, this result will provide a framework to understand the trade-off between the (rate) performance and the (quantum) receiver complexity, for reliable transmission of classical information over a quantum channel.

Theorem 5.2. *For a pure-state classical-quantum (cq) channel $W : x \rightarrow |\psi_x\rangle$, $x \in \mathcal{X}$, the maximum achievable information rate using quantum measurements of blocklength N , which is C_N/N as defined in (5.2), is lower bounded by*

$$\frac{C_N}{N} \geq \max_R \left(\left(1 - 2e^{-NE(R)}\right) R - \frac{\log 2}{N} \right), \quad (5.17)$$

where

$$E(R) = \max_{0 \leq s \leq 1} \left(\max_{P_X} (-\log \text{Tr}(\rho^{1+s})) - sR \right), \quad (5.18)$$

with $\rho = \sum_{x \in \mathcal{X}} P_X(x) |\psi_x\rangle \langle \psi_x|$.

By using this theorem, for the BPSK input channel, a blocklength N can be calculated at which the lower bound of (5.17) exceeds certain targeted rates below capacity. In the previous section, it was shown that there is a gap between $C_{1,\text{BPSK}}(\mathcal{E})/\mathcal{E}$ in (5.13) and $C_{\text{BPSK}}(\mathcal{E})/\mathcal{E}$ in (5.12) as $\mathcal{E} \rightarrow 0$:

$$\begin{aligned} \frac{C_{1,\text{BPSK}}(\mathcal{E})}{\mathcal{E}} &= 2 + o(1), \\ \frac{C_{\text{BPSK}}(\mathcal{E})}{\mathcal{E}} &= \log \frac{1}{\mathcal{E}} + 1 + o(1). \end{aligned}$$

We saw that the capacity of the BPSK alphabet is as good as that of the optimum continuous Gaussian-distributed input as N goes to infinity, i.e., $C_{\text{BPSK}}(\mathcal{E})$ is the same as $C(\mathcal{E})$ in the first two order terms. However, at the measurement blocklength $N = 1$, a BPSK constellation cannot even achieve the maximum mutual information of the opti-

mum binary input channel, $C_{1,\text{Binary}}(\mathcal{E})$ in (5.16), and the PIE caps off at 2 nats/photon. Therefore, the performance of the BPSK channel depends significantly on the regime of N . We will now find how much quantum processing is sufficient in order to communicate using the BPSK alphabet at rates close to its capacity.

The following corollary summarizes an answer for the question. Note that for the BPSK inputs $\{|\sqrt{\mathcal{E}}\rangle, |-\sqrt{\mathcal{E}}\rangle\}$, any input distribution satisfies the energy constraint of \mathcal{E} . Consequently, we can directly apply Theorem 5.2 to the BPSK channel—while automatically satisfying the energy constraint—even though the theorem itself does not assume any energy constraint.

Corollary 5.3. *For the coherent state BPSK channel with the energy constraint of \mathcal{E} where $\mathcal{E} \leq 0.01$, in the regime of $N \geq \mathcal{E}^{-1}(\log(1/\mathcal{E}))$,*

$$\frac{C_{N,\text{BPSK}}}{N} \geq \left(\left(1 - 2e^{-N\tilde{E}(R^*)}\right) R^* - \frac{\log 2}{N} \right), \quad (5.19)$$

where

$$\begin{aligned} R^* &= \mathcal{E} \log \frac{1}{\mathcal{E}} \left(1 - \sqrt{\frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}} \right) + \mathcal{E}, \\ \tilde{E}(R) &= -\log \left(\left(\frac{1+e^{-2\mathcal{E}}}{2} \right)^{1+s} + \left(\frac{1-e^{-2\mathcal{E}}}{2} \right)^{1+s} \right) - sR, \text{ for} \\ s &= \frac{\log \log(1/\mathcal{E}) - \log(R - \mathcal{E})}{\log(1/\mathcal{E})} - 1. \end{aligned}$$

Remark 5.1. *As $\mathcal{E} \rightarrow 0$, the lower bound of $C_{N,\text{BPSK}}/N$ in (5.19) can be simplified for $\mathcal{E}^{-1}(\log(1/\mathcal{E})) \leq N \leq \mathcal{E}^{-2}$ as*

$$\begin{aligned} \frac{C_{N,\text{BPSK}}}{N} &\geq \mathcal{E} \log \frac{1}{\mathcal{E}} \left(1 - \sqrt{\frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}} \right) + \mathcal{E} \\ &\quad + O \left(\frac{\mathcal{E} \log(1/\mathcal{E})}{\sqrt{N\mathcal{E} \log(N\mathcal{E})}} + \frac{\mathcal{E}}{\log(1/\mathcal{E})} \right). \end{aligned} \quad (5.20)$$

For a narrower range of N such that $\mathcal{E}^{-1}(\log(1/\mathcal{E}))^2 \leq N \leq \mathcal{E}^{-2}$, the lower bound can

be further simplified as

$$\frac{C_{N,\text{BPSK}}}{N} \geq \mathcal{E} \log \frac{1}{\mathcal{E}} \left(1 - \sqrt{\frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}} \right) + \mathcal{E} + o(\mathcal{E}). \quad (5.21)$$

From (5.20), it can be shown that at

$$N = 2\mathcal{E}^{-1} (\log(1/\mathcal{E}))^2 (\log \log(1/\mathcal{E}))^{-1}, \quad (5.22)$$

$$\frac{C_{N,\text{BPSK}}}{N} \geq \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + o\left(\mathcal{E} \log \log \frac{1}{\mathcal{E}}\right). \quad (5.23)$$

Therefore, for the above range of N , the coherent state BPSK channel can attain the PIE at least as high as $C_{1,\text{Binary}}(\mathcal{E})/\mathcal{E}$ for the first- and second-order terms, which is the maximum achievable PIE at $N = 1$ with a binary-input with mean photon number constraint \mathcal{E} .

Moreover, from (5.21), it can be shown that at

$$N = \mathcal{E}^{-1} (\log(1/\mathcal{E}))^2 (\log \log(1/\mathcal{E}))^2, \quad (5.24)$$

$$\frac{C_{N,\text{BPSK}}}{N} \geq \mathcal{E} \log \frac{1}{\mathcal{E}} + \mathcal{E} + o(\mathcal{E}). \quad (5.25)$$

Note that for the above range of N , the lower bound already approaches $C_{\text{BPSK}}(\mathcal{E})$ to the first two orders, which is the maximum information rate achievable using BPSK with an arbitrarily large length of quantum processing.

Proof. Appendix 5.B. □

Let us apply these results for the case when the average photon number transmitted per symbol, \mathcal{E} , is 0.01. Using the result of (5.19), the photon information efficiency achievable by the BPSK channel is plotted as a function of N in Fig. 5.2. For $\mathcal{E} = 0.01$, the inner product $\gamma := |\langle \sqrt{\mathcal{E}} | -\sqrt{\mathcal{E}} \rangle| = \exp[-2\mathcal{E}] = e^{-0.02}$. By plugging γ into (5.6) and (5.7), and dividing the resulting capacities by \mathcal{E} , the PIE at an arbitrarily large N is 5.55 nats/photon, and at $N = 1$, is 1.97 nats/photon. Therefore, as N increases

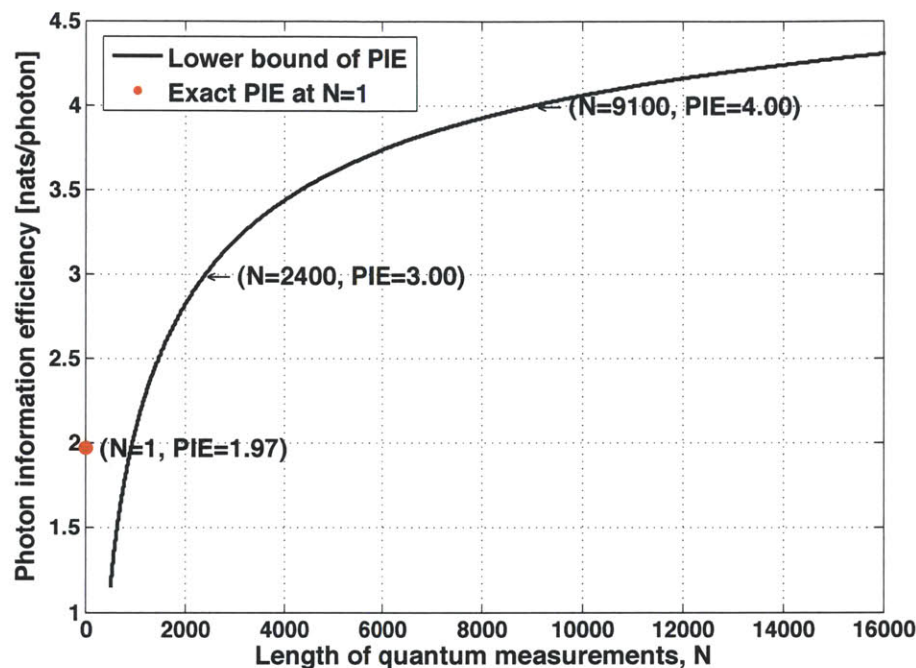


Figure 5.2. A lower bound of photon information efficiency of the BPSK channel, $C_N/(N\mathcal{E})$, at $\mathcal{E} = 0.01$ for the finite blocklength N .

from 1 to ∞ , the PIE of the BPSK channel should strictly increase from 1.97 to 5.55 nats/photon, and hence the gain in PIE from a joint measurement of an arbitrarily large length can be maximally 3.58 nats/photon. From the lower bound of PIE in Fig. 5.2, it can be seen that at $N = 2400$, a PIE of 3.0 nats/photon can be achieved, and at $N = 9100$, 4.0 nats/photon is achievable. The lower bound is not tight in the regime of very small N , but it gets tighter as N increases, and approaches the ultimate limit of PIE as $N \rightarrow \infty$.

Let us compare this exact lower bound with the approximated results from the scaling laws. From the first two order terms of the approximations in (5.13), (5.16),

and (5.12), a few reference points of PIE are calculated at $\mathcal{E} = 0.01$.

$$\begin{aligned} C_{\text{BPSK}}(\mathcal{E})/\mathcal{E} &\approx \log(1/\mathcal{E}) + 1 = 5.61 \text{ nats/photon} \\ C_{1,\text{BPSK}}(\mathcal{E})/\mathcal{E} &\approx 2.00 \text{ nats/photon} \\ C_{1,\text{Binary}}(\mathcal{E})/\mathcal{E} &\approx \log(1/\mathcal{E}) - \log \log(1/\mathcal{E}) = 3.08 \text{ nats/photon} \end{aligned} \tag{5.26}$$

We can see that these approximations of PIE for the BPSK channel at $N = 1$ and $N = \infty$ are very close to the exact calculations. Moreover, it shows that when the optimum binary input for $N = 1$ is used rather than the BPSK, the PIE of about 3.1 nats/photon is achievable even at $N = 1$. The estimated length of N to make the lower bound of $C_{N,\text{BPSK}}$ be equal to the first two order terms of $C_{1,\text{Binary}}$, C_{BPSK} , are $N = 2777$ and $N = 4946$, respectively, from (5.22) and (5.24) calculated at $\mathcal{E} = 0.01$. In Fig. 5.2, we showed that at $N = 2400$, a PIE of 3.00 nats/photon, which is close to $C_{1,\text{Binary}}(\mathcal{E})/\mathcal{E}$, is achievable. Therefore, the estimate of N from (5.22) is quite accurate at $\mathcal{E} = 0.01$. However, at $N = 4946$, the exact lower bound of PIE is still 3.67 nats/photon, which is 1.88 nats/photon away from the maximum achievable PIE with an arbitrarily large length of quantum processing. Therefore, the estimate of N in (5.24) is not very tight for \mathcal{E} in the order of 10^{-2} , but it will get tighter for smaller \mathcal{E} since it is calculated based on the assumption of $\mathcal{E} \rightarrow 0$.

■ 5.5 Proof of Theorem 5.2

Theorem 5.2 will be proved based on two lemmas that will be introduced in this section. Note that in the definition of C_N in (5.2), both the choice of the N -symbol inner code-JDR measurement pair, from which the superchannel distribution $p_{k|j}^{(N)}$ is determined, as well as the probability distribution over the inputs of the superchannel must be optimized, in order to find the maximum mutual information of the superchannel. The complexity of this optimization increases exponentially with N , and hence this optimization problem is intractable. Instead of trying to calculate the exact C_N , we provided a lower bound of C_N in the finite regime of N , which can be written as a simple

optimization over single-letter input distribution. Therefore, we can easily calculate the lower bound of C_N/N for any finite N .

The proof of Theorem 5.2 is based on two ideas: First, instead of tracking the exact superchannel distribution $p_{k|j}^{(N)}$, which depends on the detailed structure of the length- N inner code and joint measurement, we focus on one representative quantity derived from $p_{k|j}^{(N)}$ that can be easily analyzed and optimized. Second, among superchannels that have the same value of this representative quantity, we find a superchannel whose mutual information is the smallest. The representative quantity is the average decoding error probability of the inner code, under a uniform distribution over the inner codewords, defined as

$$p_e = e^{-NR} \sum_{j=1}^{e^{NR}} \sum_{k \neq j} p_{k|j}^{(N)}, \quad (5.27)$$

where R is the rate of the inner code. We first summarize previous works that investigated upper and lower bounds of p_e over N -symbol inner code and JDR measurement pairs, and then provide a lower bound on the maximum mutual information of superchannel with a fixed p_e . Theorem 5.2 will be proved by combining these two results.

■ 5.5.1 Upper and Lower Bounds on the Average Probability of Error

In ref. [24], Holevo showed the following upper bound on p_e for a code of length N and rate R .

Lemma 5.4. [Holevo] *For a pure-state classical-quantum (cq) channel $W : x \rightarrow |\psi_x\rangle$, $x \in \mathcal{X}$, there exists a block code of length N and rate R that can be decoded by a set of measurements with the average probability of error satisfying*

$$p_e \leq 2 \exp[-NE(R)], \quad (5.28)$$

where, for $\rho = \sum_x P_X(x) |\psi_x\rangle \langle \psi_x|$,

$$E(R) = \max_{0 \leq s \leq 1} \left[\max_{P_X} (-\log \text{Tr}(\rho^{1+s})) - sR \right]. \quad (5.29)$$

Proof. For reader's convenience, the proof of this lemma is summarized in Appendix 5.C. \square

Note that this result holds for *any* positive integer N . Moreover, the exponential decay rate of this upper bound on p_e is characterized by the exponent $E(R)$ that is independent of N and can be calculated from the optimization over the single-letter input distribution P_X .

Let us discuss the tightness of the upper bound in (5.28), in terms of the rate of exponential decay as $N \rightarrow \infty$. From a lower bound on p_e , it will be shown that

$$\limsup_{N \rightarrow \infty} -\frac{1}{N} \log p_e = E(R) \quad (5.30)$$

at high rates of R , i.e., $R_0 \leq R \leq C$ for a certain R_0 . For a classical discrete memoryless channel (DMC), a lower bound for the average decoding error probability of a block code was first derived by Shannon-Gallager-Berlekamp in [39], and the bound is termed *sphere packing bound*. The sphere packing bound decreases exponentially with the blocklength, and the exponent is tight at high rates below the capacity of the channel.

For quantum channels, a meaningful lower bound for p_e had not been established until very recently. In [8], a quantum analogue of the sphere packing bound was first provided based on the idea of Nussbaum-Szkola mapping, introduced in [32] as a tool to prove the converse part of the quantum Chernoff bound for binary hypothesis testing between two quantum states. The main result of [8] is summarized below.

Lemma 5.5 (Sphere packing bound for quantum channels). *When we transmit classical information over a pure-state classical-quantum (cq) channel $W : x \rightarrow |\psi_x\rangle$, $x \in \mathcal{X}$, for every length N and rate R code, the average probability of error*

$$p_e \geq \exp[-N(E_{\text{sp}}(R - \epsilon) + o(1))] \quad (5.31)$$

for every $\epsilon > 0$, where, for $\rho = \sum_x P_X(x) |\psi_x\rangle\langle\psi_x|$,

$$E_{\text{sp}}(R) = \sup_{s \geq 0} \left(\max_{P_X} (-\log \text{Tr}(\rho^{1+s})) - sR \right). \quad (5.32)$$

From the lower and upper bound of p_e in (5.28) and (5.31), we can see that when $E(R) = E_{\text{sp}}(R)$, the exponent is tight and (5.30) holds. The rates where $E(R) = E_{\text{sp}}(R)$ are in $R_0 \leq R \leq C$ where R_0 is the rate at which the optimum s to achieve $E_{\text{sp}}(R_0)$ in (5.32) equals to 1.

■ 5.5.2 Equierror Superchannel

Now, among superchannels that have the same value of p_e , we find a superchannel whose mutual information is the smallest. An *equierror superchannel*, which was first introduced in [11], is defined with the following distribution:

$$\bar{p}_{k|j}^{(N)} := \begin{cases} 1 - p_e, & k = j; \\ (e^{NR} - 1)^{-1} p_e, & k \neq j. \end{cases} \quad (5.33)$$

This channel assumes that the probability of making an error is equal for every input, and when an error occurs, all wrong estimates $k \neq j$ are equally likely. Therefore, this channel is symmetric between inputs, and is symmetric between outputs except for the right estimate, i.e., $k = j$. Due to the symmetry, the input distribution that maximizes the mutual information of this channel is uniform. The resulting maximum mutual information of the equierror superchannel,

$$\begin{aligned} \max_{p_j} I(p_j, \bar{p}_{k|j}^{(N)}) &= NR - p_e \log(e^{NR} - 1) - H_{\text{B}}(p_e) \\ &> (1 - p_e)NR - \log 2, \end{aligned} \quad (5.34)$$

where $H_{\text{B}}(p) = -p \log p - (1 - p) \log(1 - p)$.

We will now show that the mutual information of the equierror superchannel is smaller than that of any other superchannels with the same average probability of

error, p_e .

Lemma 5.6. For any $p_{k|j}^{(N)}$ with a fixed p_e defined in (5.27),

$$\max_{p_j} I(p_j, p_{k|j}^{(N)}) \geq \max_{p_j} I(p_j, \bar{p}_{k|j}^{(N)}) \quad (5.35)$$

for the equierror superchannel, $\bar{p}_{k|j}^{(N)}$ with the same p_e .

Proof. For a random variable X that is uniformly distributed over e^{NR} inputs, and the conditional distribution $P_{Y|X}(k|j) := p_{k|j}^{(N)}$,

$$\max_{p_j} I(p_j, p_{k|j}^{(N)}) \geq I(X; Y) = NR - H(X|Y). \quad (5.36)$$

From Fano's inequality, we have

$$\begin{aligned} H(X|Y) &\leq H_B(\Pr(X \neq Y)) + \Pr(X \neq Y) \log(e^{NR} - 1) \\ &= H_B(p_e) + p_e \log(e^{NR} - 1). \end{aligned}$$

By combining the above two inequalities,

$$\begin{aligned} \max_{p_j} I(p_j, p_{k|j}^{(N)}) &\geq NR - p_e \log(e^{NR} - 1) - H_B(p_e) \\ &= \max_{p_j} I(p_j, \bar{p}_{k|j}^{(N)}). \end{aligned} \quad (5.37)$$

□

Then, by the definition of C_N and Lemma 5.6, when there exists an inner code of length N and rate R that can be decoded by a set of length N measurements with an average error probability p_e ,

$$\frac{C_N}{N} \geq \max_{p_j} \frac{I(p_j, p_{k|j}^{(N)})}{N} > (1 - p_e)R - \frac{\log 2}{N}. \quad (5.38)$$

By combining Lemma 5.4 with (5.38), Theorem 5.2 can be proven.

■ 5.6 Interpretation of Superadditivity: Classical DMC vs. Quantum Channel

In Section 5.3, we demonstrated strict superadditivity of C_N , i.e., $C_N + C_M < C_{N+M}$, for binary-input quantum channels with and without an energy constraint. We provided a general lower bound on C_N/N for a fixed N in Theorem 5.2, which made it possible for us to understand the trade-off between the maximum achievable information rate and the complexity of quantum processing at the receiver as N , the length of the JDR, increases. Previously, the superadditivity of C_N has been thought of as a unique property that can be observed only in quantum channels, but not in classical DMCs. One popular interpretation of this phenomenon is that a set of length- N *entangling* quantum measurements can induce a classical superchannel that has memory over the N channel uses (despite the fact that the underlying cq channel $x \rightarrow |\psi_x\rangle$ is memoryless over each channel use). The Shannon capacity of this induced classical channel (with memory) can be higher than N times the Shannon capacity of the classical memoryless channel induced by any symbol-by-symbol receiver measurement. This capability of inducing a classical superchannel by harnessing the optimally-correlated quantum noise in the N -fold Hilbert space is what increases the number of information bits extractable per modulation symbol, when a longer block of symbols is detected collectively *while still in the quantum domain*.

Despite the fact that the above intuition of why superadditivity appears in the capacity of classical-quantum channels is somewhat satisfying, this viewpoint does not provide enough quantitative insight to fully understand the phenomenon. In this section, we will introduce a new aspect on understanding strict superadditivity of C_N by comparing the performance of concatenated coding over quantum channels as we analyze here, and concatenated coding over classical DMCs as studied by Forney [11], for a fixed inner code length N .

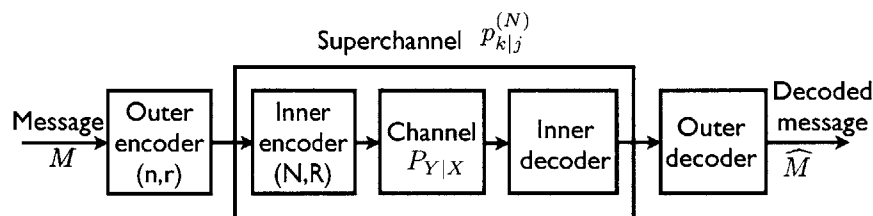


Figure 5.3. Concatenated coding over a classical DMC

■ 5.6.1 A Unifying Framework to Explain Superadditivity of C_N

Fig. 5.3 illustrates a concatenated coding architecture over a classical DMC. Compared to Fig. 5.1, the quantum channel is replaced with a classical DMC $P_{Y|X}$, and in place of a quantum joint-detection receiver, we now have a classical inner decoder. Forney introduced this concatenated coding architecture for a classical DMC to analyze the trade-off between (rate) performance and (coding) complexity for communication over a classical DMC [11]. Forney analyzed the performance by evaluating the *error exponent* achievable with a concatenated code, and also examined how the decoding complexity increases as the overall length of the concatenated code, $N_c = nN$ increases. It is obvious that when the inner decoder generates a sufficient statistic of the channel output and forwards it to the outer decoder, there is no loss of information, so that the performance of the concatenated code can be as good as an optimum code, even within the restricted structure of code concatenation. Despite the fact that the performance remains intact, the decoding complexity increases exponentially with the overall length of the code. On the other hand, it was shown in [11] that even if there is some loss of information at the inner decoder by making a hard-decision on the message of the inner code, as the inner code blocklength N goes to infinity, the capacity of the underlying classical DMC can be achieved with the concatenated code. Moreover, the overall complexity of the decoding algorithm is significantly reduced to be almost linear in the length of the concatenated code. The loss of information at the inner decoder, however, degrades the achievable error exponent over all rates below capacity.

The above result can be proved by analyzing a lower bound on the performance of

the concatenated codes over the classical DMC. To get the lower bound, the equierror superchannel defined in (5.33), whose mutual information is smaller than that of any other superchannel with the same p_e , is used. The average probability of error p_e from decoding at the inner decoder can be analyzed by using the error exponent of the classical DMC $P_{Y|X}$ in [13]. It is shown that an optimum inner code with the minimum decoding error can achieve p_e as low as

$$p_e = \exp[-N(E(R) + o(1))] \quad (5.39)$$

as $N \rightarrow \infty$, when

$$E(R) = \max_{0 \leq s \leq 1} \left(\max_{P_X} (E_0(s, P_X)) - sR \right) \quad (5.40)$$

with

$$E_0(s, P_X) := -\log \sum_y \left[\sum_x P_X(x) P_{Y|X}(y|x)^{1/(1+s)} \right]^{1+s}. \quad (5.41)$$

By using the p_e in (5.39) and analyzing the performance of the equierror channel, it can be shown that the capacity of the DMC, which is $C = \max_{P_X} I(P_X, P_{Y|X})$, can be achievable by the concatenated code *as both the inner code blocklength N and the outer code blocklength n go to infinity*, even when the inner decoder makes hard-decisions on estimating the inner code messages, and discards all the rest of the information about the channel output.

Let us clarify the difference between the concatenated code over the classical DMC and that over the quantum channel. When a likelihood detector is used at the classical inner decoder, after decoding the most likely codeword given a received channel output, the classical inner decoder can still have the information about which codeword is the second mostly likely one, and how much less likely it is compared to the first one, etc. On the contrary, for the quantum channel, once the quantum states are measured by the quantum detector, it certainly results in a loss of information since after the measurement, the quantum state of the inner codeword is lost, and in turn

all the information that was encoded in the quantum states is destroyed, except for the hard guess of the inner codeword message that the JDR generated. Therefore, different from the classical inner decoder, which has an option to maintain a sufficient statistics of the channel outputs with the cost of complexity, a loss of information at the quantum detector is not avoidable⁴. For the quantum channel, the trade-off between the achievable information rate and the complexity of quantum processing can be analyzed by observing how C_N/N increases with N . In contrast, for the concatenated code over the classical DMC, the trade-off between performance and complexity is analyzed by assuming a certain type of loss of information at the inner decoder that makes the decoding complexity increase almost linearly with the overall blocklength of the code, and by showing how the error exponent of the overall code is degraded by the loss of information at the inner decoder, *under the assumption of a large enough inner code blocklength N* .

We now ask a new question for the concatenated code over the classical DMC, similar to the one we asked for the quantum channel: When the inner decoder makes a hard estimate of messages of the inner code *at a finite blocklength N* , how does the maximum achievable information rate (error-free bits per use of the underlying DMC) with the concatenated code increase as N increases (with no restriction on the complexity of the outer code)?

For the inner decoder that makes the hard-decision at a finite blocklength N of the inner code, the maximum achievable information rate by the concatenated code is

⁴We should caveat this statement by the fact that the quantum JDR acting on the inner codeword does not *have to* generate a hard output on the inner code message. In fact, it is known that the number of outcomes in the POVM that maximizes the accessible information for M linearly independent pure states, grows as $O(M^2)$ [9]. In our case, $M = 2^{NR}$. Nonetheless, even such a POVM is a harder-decision measurement than having access to all the message likelihoods. In recent years, some quantum decoding techniques have been developed—such as the sequential decoder [15] and the quantum successive-cancellation decoder for a quantum polar code [17, 46]—that achieve the Holevo capacity, which make weak (partially-destructive) measurements on the received codeword, and retain the post-measurement states for further conditional quantum processing. Recently Wilde *et al.* used a quantum version of the likelihood ratio test, originally proposed by Fuchs and Caves [12]—another non-destructive quantum measurement—in an attempt to build an efficient decoder for the quantum polar code [48]. However, all these weak non-destructive quantum measurements are very hard to realize.

C_N/N where,

$$C_N = \max_{p_j} \max_{\{\text{N-symbol inner code-decoder pairs}\}} I(p_j, p_{k|j}^{(N)}) \quad (5.42)$$

for the superchannel distribution $p_{k|j}^{(N)}$, which is determined by the decoding algorithm, given an inner code. By using Lemma 2, it can also be shown that when there exists a code of length N and rate R whose probability of decoding error is p_e ,

$$\frac{C_N}{N} > (1 - p_e)R - \frac{\log 2}{N}. \quad (5.43)$$

Moreover, in [13], it is shown that for the classical DMC $P_{Y|X}$, there exists a code of length N and rate R whose probability of error p_e is bounded by

$$p_e \leq \exp[-NE(R)] \quad (5.44)$$

with $E(R)$ in (5.40). By combining (5.43) and (5.44), the following theorem can be proved for the maximum achievable information rate of the concatenated codes over the classical DMC at a finite N .

Theorem 5.7. *With a fixed inner code blocklength N ,*

$$\frac{C_N}{N} \geq \max_R \left(\left(1 - e^{-NE(R)}\right) R - \frac{\log 2}{N} \right), \quad (5.45)$$

with $E(R)$ as defined in (5.40).

Note that the lower bound on C_N/N in (5.45) strictly increases with N , and it has exactly the same form as that for the quantum channel in (5.17) except for the difference in $E(R)$ and a constant multiplying $e^{-NE(R)}$. As a result, we can observe a phenomenon similar to the superadditivity of C_N in the quantum channel, even in the classical DMC when the inner decoder makes hard-decisions at a finite blocklength. The reason why C_N is away from the capacity of the channel C for a finite inner code blocklength N

is because the hard-decision at the inner decoder results in a significant amount of loss of information, which even hurts the rate of the communication. Moreover, as N increases, the quality of the hard-decision is improved, which makes it possible to achieve a higher information rate. Therefore, the superadditivity of C_N can now be interpreted as a degradation of the performance by the loss of information at the inner decoder that makes the hard-decision at a finite blocklength. This new understanding can also be applied to explain the same phenomenon observed in the quantum channel by replacing the role of inner decoder with a quantum joint-detection receiver that makes hard-decisions on finite blocks of quantum states.

■ 5.6.2 An Approximation of the Lower Bound on C_N

We will simplify the lower bound of C_N by finding an approximation of the error exponent $E(R)$ for the quantum channel and for the classical DMC. Using the simplified lower bound, it will be possible to compare the quantum channel and the classical channel by calculating the inner code blocklength N required to achieve a given fraction of the ultimate capacity of each channel. To avoid confusion, from this point on, a function for the quantum channel will be written with a superscript (q) and that for the classical DMC with a superscript (c); for example, $E^{(q)}(R)$ and $E^{(c)}(R)$.

The error exponent of the classical DMC, $E^{(c)}(R)$ in (5.40), can be approximated by the Taylor expansion at the rate R close to the capacity C as

$$E^{(c)}(R) = \frac{1}{2V^{(c)}}(R - C)^2 + O((R - C)^3), \quad (5.46)$$

with a parameter $V^{(c)}$, where

$$V^{(c)} = \sum_{x,y} p_x p_{y|x} \left[\left(\log \frac{p_{y|x}}{p_y} - \sum_{x,y} p_x p_{y|x} \log \frac{p_{y|x}}{p_y} \right)^2 \right], \quad (5.47)$$

for the capacity achieving input distribution $p_x := P_X^*(x)$ and the corresponding capacity achieving output distribution $p_y := P_Y^*(y)$ according to the channel $p_{y|x} :=$

$P_{Y|X}(y|x)$. In (5.47), $V^{(c)}$ is the variance of $\log(p_{y|x}/p_y)$ under the distribution $p_x p_{y|x}$, and was termed the *channel dispersion* in [34].

Similarly, the error exponent of the quantum channel, $E^{(q)}(R)$ in (5.18), can be approximated with a parameter $V^{(q)}$, which is a characteristic of the quantum channel similar to the channel dispersion of the classical DMC. The definition of $V^{(q)}$ depends on the average density operator ρ , which fully characterizes the classical capacity of the pure-state quantum channel. For a set of input states $\{|\psi_x\rangle\}$, when P_X^* is the optimum input distribution that attains the capacity of the quantum channel $C = \max_{P_X} \text{Tr}(-\rho \log \rho)$ where $\rho = \sum_x P_X(x) |\psi_x\rangle \langle \psi_x|$, the parameter $V^{(q)}$ is defined by the eigenvalues of the density operator ρ at $P_X = P_X^*$. Let us denote the eigenvalues of ρ by σ_i , $i = 1, \dots, J$ where J is the dimension of the space spanned by the input states $\{|\psi_x\rangle\}$. From the fact that ρ is a positive operator and $\text{Tr}(\rho) = 1$, it can be shown that each $\sigma_i \geq 0$ for all i and $\sum_{i=1}^J \sigma_i = 1$. Then, $V^{(q)}$ is defined as a variance of the random variable $-\log \sigma$ where $\sigma \in \{\sigma_i\}$ with probability distribution $\{\sigma_1, \dots, \sigma_J\}$, i.e.,

$$V^{(q)} = \sum_{i=1}^J \sigma_i (-\log \sigma_i)^2 - \left(\sum_{i=1}^J \sigma_i (-\log \sigma_i) \right)^2. \quad (5.48)$$

By the Taylor expansion of $E^{(q)}(R)$ in (5.18) at the rate R close to C , it can be shown that

$$E^{(q)}(R) = \frac{1}{2V^{(q)}}(R - C)^2 + O((R - C)^3). \quad (5.49)$$

Therefore, both the error exponent of the classical DMC and that of the quantum channel can be approximated as a quadratic term in the rate R with the quadratic coefficient inversely proportional to the dispersion of the channel. Since the lower bound on C_N as well as the approximated error exponent $E(R)$ have similar forms for the classical DMC and for the quantum channel, it is possible to compare the classical DMC and the quantum channel by a common simplified lower bound on C_N , which can be written with the parameter V and C as follows.

Theorem 5.8. *For both a classical DMC and a pure-state classical-quantum channel,*

when the channel dispersion V and the capacity C satisfy i) $\sqrt{\frac{V}{NC^2}} \rightarrow 0$ as $N \rightarrow \infty$ and ii) $V \cdot C$ is finite, the maximum achievable information rate at the inner code blocklength N is lower bounded by

$$\frac{C_N}{N} \geq C \cdot \left(1 - \sqrt{\frac{V}{NC^2} \log \left(\frac{NC^2}{V} \right)} \right) - \frac{\log 2}{N} + O \left(\sqrt{\frac{V}{N \log \frac{NC^2}{V}}} \log \log \left(\frac{NC^2}{V} \right) \right). \quad (5.50)$$

Proof. The quadratic approximation of $E(R)$ can be used to find a simplified form for a lower bound on C_N/N . Both for the quantum channel and the classical channel, C_N is lower bounded by

$$\frac{C_N}{N} \geq \max_R \left(\left(1 - 2e^{-NE(R)} \right) R - \frac{\log 2}{N} \right), \quad (5.51)$$

using Theorems 5.2 and 5.7. Then, for a fixed rate

$$R^* = C \cdot \left(1 - \sqrt{\frac{V}{NC^2} \log \left(\frac{NC^2}{V} \log \frac{NC^2}{V} \right)} \right), \quad (5.52)$$

whose derivation is omitted in this chapter, the approximated error exponent at R^* is

$$\begin{aligned} E(R^*) &= \frac{1}{2N} \log \left(\frac{NC^2}{V} \log \frac{NC^2}{V} \right) \\ &+ O \left(\frac{VC}{N} \sqrt{\frac{V}{NC^2}} \left(\log \left(\frac{NC^2}{V} \log \frac{NC^2}{V} \right) \right)^{3/2} \right). \end{aligned} \quad (5.53)$$

It can be checked that under the assumptions of i) $\sqrt{\frac{V}{NC^2}} \rightarrow 0$ and ii) $V \cdot C$ is finite, the term in $O(\cdot)$ in (5.53) approaches 0, which results in

$$e^{-NE(R^*)} = \sqrt{\frac{V}{NC^2 \log(NC^2/V)}} (1 + o(1)). \quad (5.54)$$

By plugging (5.52) and (5.54) into the lower bound (5.51), C_N/N can be bounded

as shown in (5.50). \square

Remark 5.2. *From the lower bound of Theorem 5.8, we can see that the inner code blocklength N at which the lower bound is equal to a given fraction of the capacity is proportional to V/C^2 .*

Since the same bound on C_N/N as in (5.50) holds both for the quantum and the classical channels, using the parameter V/C^2 , we can compare the behaviors of the quantum channel and the classical DMC. For the BPSK quantum channel, by using the two eigenvalues of ρ at P_X^* , which are $\sigma_1 = (1 - e^{-2\mathcal{E}})/2$ and $\sigma_2 = (1 + e^{-2\mathcal{E}})/2$, the channel dispersion and the capacity can be calculated as

$$\begin{aligned} V_{\text{BPSK}} &= \mathcal{E} \left(\log \frac{1}{\mathcal{E}} \right)^2 (1 + O(\mathcal{E})), \text{ and} \\ C_{\text{BPSK}} &= \mathcal{E} \log \frac{1}{\mathcal{E}} + \mathcal{E} + o(\mathcal{E}). \end{aligned} \tag{5.55}$$

Then, $V_{\text{BPSK}}/C_{\text{BPSK}}^2 \approx 1/\mathcal{E}$ for the low photon number regime where $\mathcal{E} \rightarrow 0$. For the classical additive white Gaussian noise (AWGN) channel in the low-power regime where $\text{SNR} \rightarrow 0$, $V_{\text{AWGN}}/C_{\text{AWGN}}^2$ can be calculated by using the result of [34], and it is $4/\text{SNR}$. For both channels, V/C^2 is inversely proportional to the energy to transmit the information per channel use. This means that as the energy decreases, in order to make the lower bound meet a targeted fraction of capacity, it is necessary to adopt a longer inner code.

■ 5.7 Conclusion

The Holevo capacity of a classical-quantum channel, i.e., the ultimate rate of reliable communication for sending classical data over a quantum channel, is a doubly-asymptotic result; meaning the achievability of the capacity C has been proven so far for the case when the transmitter is allowed to code over an arbitrarily large sequence of quantum states (spanning N_c channel uses), *and* when the receiver is assumed to be able to *jointly* measure quantum states of the received codewords, also over N_c

channel uses, while $N_c \rightarrow \infty$. However, the assumption that arbitrarily large number of quantum states can be jointly measured is the primary barrier prohibiting practical implementations of joint detection receivers—particularly in the context of optical communication. Our goal in this chapter was to separate these two infinities: the coding blocklength N_c (a relatively inexpensive resource), and the length of the joint detection receiver, $N \leq N_c$ (a far more expensive resource), and to evaluate how the capacity C_N , constrained to length- N joint measurements (but no restrictions on the classical code complexity), grows with N . We analyzed superadditivity in classical capacity of a pure-state quantum channel while focusing on the quantitative trade-off between reliable-rate performance and quantum-decoding complexity. In order to analyze this trade-off, we adopted a concatenated coding scheme where a quantum joint-detection receiver acts on finite-blocklength quantum codewords of the inner code, and we found a lower bound on the maximum achievable information rate as a function of the length N of the quantum measurement that decodes the inner code. We also observed a similar phenomenon for a classical discrete memoryless channel (DMC), and explained how a classical superadditivity in channel capacity occurs due to a loss of information from the hard-decision at the inner decoder of finite blocklength N . We developed a unifying framework, within which the superadditivity in capacity of the classical DMC and that of the pure-state quantum channel can be compared with a parameter V/C^2 (where V is the channel dispersion, and C is channel capacity), which is proportional to the inner-code measurement N that is sufficient to achieve a given fraction of the capacity.

■ 5.A Proof of Lemma 5.1

$C_{1,\text{Binary}}(\mathcal{E})$ is the maximum mutual information of binary input channels $\{|\alpha_0\rangle, |\alpha_1\rangle\}$ under the average photon number constraint, $(1-q)|\alpha_0|^2 + q|\alpha_1|^2 \leq \mathcal{E}$. In [41], it is shown that

$$C_{1,\text{Binary}}(\mathcal{E}) = \max_{(1-q)|\alpha_0|^2 + q|\alpha_1|^2 = \mathcal{E}} H_B(q) - H_B(p) \quad (5.56)$$

where $H_B(x) = -x \log x - (1-x) \log(1-x)$, and

$$p = \frac{1 - \sqrt{1 - 4q(1-q)e^{-|\alpha_0 - \alpha_1|^2}}}{2}. \quad (5.57)$$

We first find the optimum input states $\{|\alpha_0\rangle, |\alpha_1\rangle\}$ for a fixed q and then will find q that maximizes (5.56). In [41], it was shown how to find the optimum input states given a fixed q , but we summarize the steps again for reader's convenience.

For a fixed q , the input states that maximize $C_{1,\text{Binary}}(\mathcal{E})$ should minimize $H_B(p)$. Since $H_B(p)$ is an increasing function in p for $0 \leq p \leq 1/2$, and p in (5.57) gets smaller as $|\alpha_0 - \alpha_1|^2$ increases, we need find inputs $\alpha_0, \alpha_1 \in \mathbb{C}$ that maximizes $|\alpha_0 - \alpha_1|^2$ under the energy constraint $(1-q)|\alpha_0|^2 + q|\alpha_1|^2 \leq \mathcal{E}$ for a given q . To maximize $|\alpha_0 - \alpha_1|^2$ under the photon number constraint,

$$\alpha_1 = -k\alpha_0 \quad (5.58)$$

for a real number $k \geq 0$ that satisfies

$$(1-q)|\alpha_0|^2 + q|\alpha_1|^2 = (1-q+k^2 \cdot q)|\alpha_0|^2 = \mathcal{E}. \quad (5.59)$$

The optimum k that maximizes $f(k) := |\alpha_0 - \alpha_1|^2 = (1+k)^2|\alpha_0|^2 = ((1+k)^2\mathcal{E}) / (1-q+k^2 \cdot q)$ can be found from

$$\frac{\partial f(k)}{\partial k} = 0, \quad (5.60)$$

and the solution is $k^* = (1-q)/q$. Therefore, the optimum inputs and the resulting p in (5.57) become

$$\begin{aligned} \alpha_0^* &= \sqrt{\mathcal{E} \cdot q / (1-q)} \\ \alpha_1^* &= -\sqrt{\mathcal{E} \cdot (1-q) / q} \\ p^* &= \left(1 - \sqrt{1 - 4q(1-q) \exp\left(-\frac{\mathcal{E}}{q(1-q)}\right)} \right) / 2 \end{aligned} \quad (5.61)$$

Now, $C_{1,\text{Binary}}(\mathcal{E})$ can be written as

$$C_{1,\text{Binary}}(\mathcal{E}) = \max_q H_B(q) - H_B(p^*) \quad (5.62)$$

for $0 \leq q \leq 1/2$.

Define $I(q) := H_B(q) - H_B(p^*)$. The optimum q^* maximizing $I(q)$ should satisfy $\partial I(q)/\partial q|_{q=q^*} = 0$. However,

$$\begin{aligned} \frac{\partial I(q)}{\partial q} = & \log \frac{1-q}{q} - \log \frac{1-p^*}{p^*} \times \\ & \left(\frac{1-2q}{1-2p^*} \left(1 + \frac{\mathcal{E}}{q(1-q)} \right) \exp \left(-\frac{\mathcal{E}}{q(1-q)} \right) \right), \end{aligned} \quad (5.63)$$

and $\partial I(q)/\partial q = 0$ does not have a closed form solution. Instead, by focusing on the low photon number regime where $\mathcal{E} \rightarrow 0$, we can find an approximate solution for q^* , and calculate $C_{1,\text{Binary}}(\mathcal{E})$ for the first and second order terms. First, it will be shown that

$$\left. \frac{\partial I(q)}{\partial q} \right|_{q=\frac{\mathcal{E}}{2}\sqrt{\log \frac{1}{\mathcal{E}}}} > 0 \quad \text{and} \quad \left. \frac{\partial I(q)}{\partial q} \right|_{q=\frac{\mathcal{E}}{2}(\log \frac{1}{\mathcal{E}})} < 0, \quad (5.64)$$

which implies that the optimum q^* is between

$$\frac{\mathcal{E}}{2}\sqrt{\log \frac{1}{\mathcal{E}}} \leq q^* \leq \frac{\mathcal{E}}{2} \left(\log \frac{1}{\mathcal{E}} \right) \quad (5.65)$$

Let us write down the approximation of $\partial I(q)/\partial q$ as $\mathcal{E} \rightarrow 0$. When $\frac{\mathcal{E}}{2}\sqrt{\log \frac{1}{\mathcal{E}}} \leq q \leq \frac{\mathcal{E}}{2}(\log \frac{1}{\mathcal{E}})$, by Taylor expansion, each term in (5.63) can be approximated as

$$\begin{aligned} \exp \left(-\frac{\mathcal{E}}{q(1-q)} \right) &= 1 - \mathcal{E}/q + \mathcal{E}^2/(2q^2) + O(\mathcal{E}^3/q^3), \\ p^* &= q \left(1 - \mathcal{E}/q + \mathcal{E}^2/(2q^2) \right) + O(\mathcal{E}^3/q^2), \\ \log \frac{1-q}{q} &= \log(1/q) - q + O(q^2), \\ \log \frac{1-p^*}{p^*} &= \log(1/q) + \mathcal{E}/q - \mathcal{E}^2/(2q^2) + O(\mathcal{E}^3/q^3), \end{aligned}$$

$$\begin{aligned}\frac{1-2q}{1-2p^*} &= 1 + O(\mathcal{E}), \\ 1 + \frac{\mathcal{E}}{q(1-q)} &= 1 + \mathcal{E}/q + O(\mathcal{E}).\end{aligned}$$

By using these approximations, it can be shown that $\partial I(q)/\partial q$ in (5.63) can be written as

$$\frac{\partial I(q)}{\partial q} = \frac{\mathcal{E}^2}{2q^2} \log \frac{1}{q} - \frac{\mathcal{E}}{q} + O\left(\frac{\mathcal{E}^3}{q^3} \log \frac{1}{q}\right). \quad (5.66)$$

Now, it can be checked that when $q = (\mathcal{E} \sqrt{\log(1/\mathcal{E})})/2$,

$$\begin{aligned}\frac{\partial I(q)}{\partial q} &= \frac{2}{\log \frac{1}{\mathcal{E}}} \left(\log \frac{2}{\mathcal{E}} - \frac{1}{2} \log \log \frac{1}{\mathcal{E}} \right) + O\left(\frac{1}{\sqrt{\log \frac{1}{\mathcal{E}}}}\right), \\ &= 2 + o(1) > 0,\end{aligned} \quad (5.67)$$

and when $q = (\mathcal{E} \log(1/\mathcal{E}))/2$,

$$\begin{aligned}\frac{\partial I(q)}{\partial q} &= \frac{2}{(\log \frac{1}{\mathcal{E}})^2} \left(\log \frac{2}{\mathcal{E}} - \log \log \frac{1}{\mathcal{E}} \right) - \frac{2}{\log \frac{1}{\mathcal{E}}} \\ &\quad + O\left(\frac{1}{(\log \frac{1}{\mathcal{E}})^2}\right) \\ &= -\frac{2 \log \log \frac{1}{\mathcal{E}}}{(\log \frac{1}{\mathcal{E}})^2} + O\left(\frac{1}{(\log \frac{1}{\mathcal{E}})^2}\right) < 0.\end{aligned} \quad (5.68)$$

Therefore, (5.65) is verified.

Now, we will write $q^* = \frac{\mathcal{E}}{2} (\log \frac{1}{\mathcal{E}})^\alpha$ for some $1/2 \leq \alpha \leq 1$, and then will find α^* that maximizes $I(q)$. Let us find the approximation of $I(q)$ as $\mathcal{E} \rightarrow 0$. By using the Taylor expansion for $H_B(x) = -x \log x + x + O(x^2)$ as $x \rightarrow 0$,

$$H_B(p^*) = -q \log q + \mathcal{E} \log q - (\mathcal{E}^2 \log q)/(2q) + q + O(\mathcal{E}), \quad (5.69)$$

and thus

$$\begin{aligned} I(q) &= H_B(q) - H_B(p^*) \\ &= -\mathcal{E} \log q + (\mathcal{E}^2 \log q)/(2q) + O(\mathcal{E}). \end{aligned} \quad (5.70)$$

At $q = q^* = \frac{\mathcal{E}}{2} (\log \frac{1}{\mathcal{E}})^\alpha$,

$$I(q) = \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \left(\log \frac{1}{\mathcal{E}} \right)^{1-\alpha} - \alpha \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}). \quad (5.71)$$

It can be easily checked that $\partial I(q)/\partial \alpha = 0$ when $\alpha = 1$ from

$$\partial I(q)/\partial \alpha = \mathcal{E} \log \log \frac{1}{\mathcal{E}} \left(1 - \left(\log \frac{1}{\mathcal{E}} \right)^{1-\alpha} \right). \quad (5.72)$$

Therefore, the optimum $q^* = \frac{\mathcal{E}}{2} (\log \frac{1}{\mathcal{E}})$ and by plugging $\alpha = 1$ in (5.71),

$$\begin{aligned} C_{1,\text{Binary}} &= \max_q I(q) = I(q)|_{\alpha=1} \\ &= \mathcal{E} \log \frac{1}{\mathcal{E}} - \mathcal{E} \log \log \frac{1}{\mathcal{E}} + O(\mathcal{E}). \end{aligned} \quad (5.73)$$

■ 5.B Proof of Corollary 5.3

For the BPSK inputs $\{|\sqrt{\mathcal{E}}\rangle, |-\sqrt{\mathcal{E}}\rangle\}$ with input distribution $\{1-q, q\}$, two eigenvalues of the density operator $\rho = (1-q)|\sqrt{\mathcal{E}}\rangle\langle\sqrt{\mathcal{E}}| + q|-\sqrt{\mathcal{E}}\rangle\langle-\sqrt{\mathcal{E}}|$ are

$$\begin{aligned} \sigma_1 &= \left(1 - \sqrt{1 - 4q(1-q)(1 - e^{-4\mathcal{E}})} \right) / 2, \\ \sigma_2 &= \left(1 + \sqrt{1 - 4q(1-q)(1 - e^{-4\mathcal{E}})} \right) / 2, \end{aligned} \quad (5.74)$$

from (5.4) and $\gamma = |\langle\sqrt{\mathcal{E}}|-\sqrt{\mathcal{E}}\rangle| = e^{-2\mathcal{E}}$.

It can be easily checked that the optimum q that maximizes

$$-\log \text{Tr}(\rho^{1+s}) = -\log (\sigma_1^{1+s} + \sigma_2^{1+s}) \quad (5.75)$$

is equal to $1/2$ by using the symmetry between q and $1 - q$ in (5.74).

When σ_1 and σ_2 at $q = 1/2$ are denoted as σ_1^* and σ_2^* ,

$$\begin{aligned}\sigma_1^* &= (1 - e^{-2\mathcal{E}})/2, \\ \sigma_2^* &= (1 + e^{-2\mathcal{E}})/2.\end{aligned}\tag{5.76}$$

Then, the error exponent $E(R)$ in (5.18) for the BPSK inputs can be written as

$$\begin{aligned}E(R) &= \max_{0 \leq s \leq 1} \left(\max_{P_X} (-\log \text{Tr}(\rho^{1+s})) - sR \right) \\ &= \max_{0 \leq s \leq 1} \left(-\log ((\sigma_1^*)^{1+s} + (\sigma_2^*)^{1+s}) - sR \right)\end{aligned}\tag{5.77}$$

A closed form solution for the optimum s that achieves $E(R)$ cannot be found. Instead, by using the assumption of the low photon number regime, i.e., $\mathcal{E} \ll 1$, we pick the following s' and find a lower bound of $E(R)$.

$$s' := \begin{cases} \frac{\log \log(1/\mathcal{E}) - \log(R - \mathcal{E})}{\log(1/\mathcal{E})} - 1, & R_c < R < C; \\ 1, & R \leq R_c; \\ 0, & R \geq C \end{cases}\tag{5.78}$$

where $R_c = \mathcal{E} + \mathcal{E}^2 \log(1/\mathcal{E})$ and $C = \mathcal{E} \log(1/\mathcal{E}) + \mathcal{E}$. When we define

$$\tilde{E}(R) := \left(-\log \left((\sigma_1^*)^{1+s'} + (\sigma_2^*)^{1+s'} \right) - s'R \right),\tag{5.79}$$

the error exponent $E(R)$ is lower bounded by $\tilde{E}(R)$ for every R , i.e.,

$$E(R) \geq \tilde{E}(R).\tag{5.80}$$

Therefore, the lower bound of C_N/N in (5.17) from Theorem 5.2 can be further lower

bounded by using $\tilde{E}(R)$ as follows.

$$\begin{aligned} \frac{C_N}{N} &\geq \max_R \left((1 - 2e^{-NE(R)})R - \frac{\log 2}{N} \right) \\ &\geq \max_R \left((1 - 2e^{-N\tilde{E}(R)})R - \frac{\log 2}{N} \right) \end{aligned} \quad (5.81)$$

A closed form solution of the optimum R that maximizes the lower bound in (5.81) cannot be found, but again by using the assumption that $\mathcal{E} \ll 1$, we pick

$$R^* = \mathcal{E} \log \frac{1}{\mathcal{E}} \left(1 - \sqrt{\frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}} \right) + \mathcal{E} \quad (5.82)$$

for $N \geq \mathcal{E}^{-1} \log(1/\mathcal{E})$. It can be shown that for $\mathcal{E} \leq e^{-2} \approx 0.13$, the chosen rate R^* is in $R_c < R^* < C$ where $N \geq \mathcal{E}^{-1} \log(1/\mathcal{E})$, and thus s' at $R = R^*$ belongs to the first case in (5.78). To show this, we use the fact that for $N\mathcal{E} \geq 2$,

$$0 \leq \sqrt{\frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}} \leq 0.85, \quad (5.83)$$

which can be validated by numerical calculations using a computer. Since we assume that $N \geq \mathcal{E}^{-1} \log(1/\mathcal{E})$, if $\log(1/\mathcal{E}) \geq 2$, i.e., $\mathcal{E} \leq e^{-2}$, then $N\mathcal{E} \geq 2$. Moreover, if

$$\sqrt{\frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}} < 1 - \mathcal{E}, \quad (5.84)$$

then R^* in (5.82) stays in the range of $R_c < R^* < C$. From (5.83), the inequality in (5.83) holds for $\mathcal{E} < 0.15$.

Therefore, for $\mathcal{E} \leq \min\{0.15, e^{-2}\} = e^{-2}$,

$$\frac{C_N}{N} \geq (1 - 2e^{-N\tilde{E}(R^*)})R^* - \frac{\log 2}{N} \quad (5.85)$$

in the range of $N \geq \mathcal{E}^{-1} \log(1/\mathcal{E})$. By numerical calculations, we checked that the lower bound of (5.85) strictly increases with N if $\mathcal{E} \leq 0.01$. Even though the lower bound

itself is valid for $\mathcal{E} \leq e^{-2}$, since R^* and s' were chosen based on the assumption that $\mathcal{E} \ll 1$, the lower bound becomes meaningful in the sense that it strictly increases with N for small range of \mathcal{E} . Therefore, we state the corollary with the assumption that $\mathcal{E} \leq 0.01$. The corollary is proven.

Now, we will show the approximation of (5.85) as $\mathcal{E} \rightarrow 0$. For $0 < s < 1$, by using Taylor expansion,

$$\begin{aligned} (\sigma_1^*)^{1+s} &= 1 - (1+s)\mathcal{E} + \frac{(1+s)(2+s)}{2}\mathcal{E}^2 + O(\mathcal{E}^3), \\ (\sigma_2^*)^{1+s} &= \mathcal{E}^{1+s} - (1+s)\mathcal{E}^{2+s} + O(\mathcal{E}^{3+s}). \end{aligned} \quad (5.86)$$

By using these approximations and Taylor expansion of $\log(1+x) = x + O(x^2)$ as $x \rightarrow 0$,

$$-\log((\sigma_1^*)^{1+s} + (\sigma_2^*)^{1+s}) = (1+s)\mathcal{E} - \mathcal{E}^{1+s} + O(\mathcal{E}^2). \quad (5.87)$$

Then, for $s = s'$ in the range of $R_c < R < C$,

$$\begin{aligned} \tilde{E}(R) &= (1+s')\mathcal{E} - \mathcal{E}^{1+s'} - s'R + O(\mathcal{E}^2) \\ &= \frac{(R-\mathcal{E})}{\log(1/\mathcal{E})} \left(\log(R-\mathcal{E}) + \log \frac{1}{\mathcal{E}} - \log \log \frac{1}{\mathcal{E}} - 1 \right) + \mathcal{E} + O(\mathcal{E}^2). \end{aligned} \quad (5.88)$$

Now, at $R = R^*$,

$$\tilde{E}(R^*) = \mathcal{E} \cdot \left(\sqrt{f} + \log(1-\sqrt{f}) - \sqrt{f} \log(1-\sqrt{f}) \right) + O(\mathcal{E}^2) \quad (5.89)$$

where

$$f = \frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}. \quad (5.90)$$

In the range of $N \geq \mathcal{E}^{-1} \log(1/\mathcal{E})$, i.e., $N\mathcal{E} \geq \log(1/\mathcal{E})$, the resulting $f \rightarrow 0$ as $\mathcal{E} \rightarrow 0$, and thus it can be approximated as

$$\tilde{E}(R^*) = (\mathcal{E} \cdot f)/2 + O(\mathcal{E} \cdot f^{3/2} + \mathcal{E}^2). \quad (5.91)$$

If we further restrict the range of N such that

$$\mathcal{E}^{-1} \log(1/\mathcal{E}) \leq N \leq \mathcal{E}^{-2}, \text{ i.e., } \log(1/\mathcal{E}) \leq N\mathcal{E} \leq \mathcal{E}^{-1},$$

$\tilde{E}(R^*)$ becomes

$$\tilde{E}(R^*) = \frac{\mathcal{E}}{2} \cdot \frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}} + O\left(\frac{1}{N}\right) \quad (5.92)$$

as $\mathcal{E} \rightarrow 0$. Therefore,

$$\begin{aligned} N\tilde{E}(R^*) &= \log \sqrt{N\mathcal{E} \log(N\mathcal{E})} + O(1) \\ e^{-N\tilde{E}(R^*)} &= O\left(\frac{1}{\sqrt{N\mathcal{E} \log(N\mathcal{E})}}\right). \end{aligned} \quad (5.93)$$

By using this result, the lower bound of C_N/N in (5.85) can be simplified as

$$\frac{C_N}{N} \geq \mathcal{E} \log \frac{1}{\mathcal{E}} \left(1 - \sqrt{\frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}}\right) + \mathcal{E} + O\left(\frac{\mathcal{E} \log(1/\mathcal{E})}{\sqrt{N\mathcal{E} \log(N\mathcal{E})}} + \frac{\mathcal{E}}{\log(1/\mathcal{E})}\right) \quad (5.94)$$

in the range of N such that $\mathcal{E}^{-1} \log(1/\mathcal{E}) \leq N \leq \mathcal{E}^{-2}$. Moreover, for a narrower regime of N where $\mathcal{E}^{-1}(\log(1/\mathcal{E}))^2 \leq N \leq \mathcal{E}^{-2}$, the term $O\left(\frac{\mathcal{E} \log(1/\mathcal{E})}{\sqrt{N\mathcal{E} \log(N\mathcal{E})}} + \frac{\mathcal{E}}{\log(1/\mathcal{E})}\right)$ can be simplified as $o(\mathcal{E})$, and thus

$$\frac{C_N}{N} \geq \mathcal{E} \log \frac{1}{\mathcal{E}} \left(1 - \sqrt{\frac{\log(N\mathcal{E} \log(N\mathcal{E}))}{N\mathcal{E}}}\right) + \mathcal{E} + o(\mathcal{E}). \quad (5.95)$$

■ 5.C Proof of Lemma 5.4

Let us first introduce some notations related to the quantum codewords of length- N and rate- R codes. The encoder $f: \{1, \dots, M := 2^{NR}\} \rightarrow \mathcal{X}^N$ maps each message into a length- N codeword. The codeword for the j -th message can be written as $f(j) = (x_1(j), \dots, x_N(j))$ where $x_i(j) \in \mathcal{X}$ for $i = 1, \dots, N$. After sending each coded symbol through the classical-quantum channel $W: x \rightarrow |\psi_x\rangle$, the received length- N sequence

of states can be written as a density operator form,

$$S_{f(j)} := |\psi_{x_1(j)}\rangle\langle\psi_{x_1(j)}| \otimes \cdots \otimes |\psi_{x_N(j)}\rangle\langle\psi_{x_N(j)}|. \quad (5.96)$$

When we denote $|\psi_{f(j)}\rangle := |\psi_{x_1(j)}\rangle \otimes \cdots \otimes |\psi_{x_N(j)}\rangle$,

$$S_{f(j)} = |\psi_{f(j)}\rangle\langle\psi_{f(j)}|. \quad (5.97)$$

Let us define a matrix Ψ such that its j -th column is $|\psi_{f(j)}\rangle$. When $|\psi_{f(j)}\rangle$, $j \in \{1, \dots, M\}$ stay in a d -dimensional Hilbert space, the singular value decomposition of Ψ can be represented as

$$\Psi = (|\psi_{f(1)}\rangle, |\psi_{f(2)}\rangle, \dots, |\psi_{f(M)}\rangle) = U\Sigma V^\dagger, \quad (5.98)$$

where U and V are unitary matrices of size $d \times d$ and $M \times M$, respectively, and Σ is a $d \times M$ rectangular diagonal matrix with non-negative real numbers on the diagonal. V^\dagger is the Hermitian conjugate of V . We also denote the Gram matrix, Γ , and the Gram operator, G , of Ψ as

$$\begin{aligned} \Gamma &= \Psi\Psi^\dagger = V(\Sigma^\dagger\Sigma)V^\dagger, \\ G &= \Psi^\dagger\Psi = U(\Sigma\Sigma^\dagger)U^\dagger. \end{aligned} \quad (5.99)$$

Note that Γ and G are positive operators.

Now, we will introduce Square Root Measurements (SRM) $\{\Pi_j\}$, with which the M encoded quantum states are measured. The SRM is defined as a rank-one operator such that

$$\Pi_j = |\omega_j\rangle\langle\omega_j| = \left(G^{-1/2}\right) S_{f(j)} \left(G^{-1/2}\right) = \left(G^{-1/2}|\psi_{f(j)}\rangle\right) \left(\langle\psi_{f(j)}|G^{-1/2}\right) \quad (5.100)$$

where

$$G^{-1/2} = U \left((\Sigma \Sigma^\dagger)^{-1/2} \right) U^\dagger, \quad (5.101)$$

and $(\Sigma \Sigma^\dagger)^{-1/2}$ is formed by replacing every non-zero diagonal entry of $\Sigma \Sigma^\dagger$ with one over the square root of each entry. Note that the defined measurement $\{\Pi_j\}$ satisfies $\Pi_j \geq 0$, for all $j \in \{1, \dots, M\}$, and $\sum_{j=1}^M \Pi_j = \mathbb{1}$.

Define a matrix Ω such that its j -th column is the j -th measurement vector $|\omega_j\rangle$,

$$\Omega = (|\omega_1\rangle, |\omega_2\rangle, \dots, |\omega_M\rangle) = G^{-1/2} \Psi = U (\Sigma \Sigma^\dagger)^{-1/2} \Sigma V^\dagger = U \Sigma (\Sigma^\dagger \Sigma)^{-1/2} V^\dagger. \quad (5.102)$$

Then, (k, j) -entry of $\Omega^\dagger \Psi$ becomes $\langle \omega_k | \psi_{f(j)} \rangle$ and

$$\Omega^\dagger \Psi = V (\Sigma^\dagger \Sigma)^{-1/2} (\Sigma^\dagger \Sigma) V^\dagger = V (\Sigma^\dagger \Sigma)^{1/2} V^\dagger = \Gamma^{1/2}. \quad (5.103)$$

Therefore, the probability that the decoder chooses the k -th message by the SRM, when the j -th message is the true one, which is denoted as $p_{k|j}^{(N)}$, is

$$p_{k|j}^{(N)} = |\langle \omega_k | \psi_{f(j)} \rangle|^2 = |\sqrt{\Gamma_{k,j}}|^2 \quad (5.104)$$

where $\sqrt{\Gamma_{k,j}}$ denotes the (k, j) -entry of $\Gamma^{1/2}$. It means that under the SRM, once we have the geometric structure of the encoded quantum states, which is represented by its Gram matrix Γ , the distribution of the measurement outputs, given the true message, can be directly calculated from $\Gamma^{1/2}$.

From (5.104), the average probability of decoding error becomes,

$$P_e = \frac{1}{M} \sum_{j=1}^M (1 - |\sqrt{\Gamma_{j,j}}|^2) = \frac{1}{M} \sum_{j=1}^M (1 - \sqrt{\Gamma_{j,j}})(1 + \sqrt{\Gamma_{j,j}}) \quad (5.105)$$

Since $\sqrt{\Gamma_{j,j}} = \langle \omega_j | \psi_{f(j)} \rangle = \langle \psi_{f(j)} | G^{-1/2} | \psi_{f(j)} \rangle$ with the positive operator $G^{-1/2}$, and $\sum_{k=1}^M |\sqrt{\Gamma_{k,j}}|^2 = 1$, the resulting $\sqrt{\Gamma_{j,j}}$, $j = 1, \dots, M$, are positive numbers in $[0, 1]$

We will show an existence of a length- N and rate- R code of which the average

probability of error under the SRM satisfies

$$P_e \leq 2e^{-NE(R)}, \quad (5.106)$$

where

$$E(R) = \max_{0 \leq s \leq 1} \left(\max_{P_X} (-\log \text{Tr}(\rho^{1+s})) - sR \right) \quad (5.107)$$

for $\rho = \sum_x P_X(x) |\psi_x\rangle\langle\psi_x|$ as stated in Lemma 5.4.

We can further bound P_e in (5.105) as

$$P_e \leq \frac{2}{M} \sum_{j=1}^M (1 - \sqrt{\Gamma_{j,j}}) = \frac{2}{M} \left(M - \text{Tr}(\Gamma^{1/2}) \right) = \frac{2}{M} \left(M - \text{Tr}(G^{1/2}) \right). \quad (5.108)$$

Note that all the eigenvalues of $G^{1/2}$ are positive. Therefore, using the following inequality,

$$2\sqrt{x} \geq 2x - \min\{(x^2 - x), 2x\}, \quad x \geq 0, \quad (5.109)$$

we have

$$-2\text{Tr}(G^{1/2}) \leq -2\text{Tr}(G) + \min\{\text{Tr}(G^2 - G), 2\text{Tr}(G)\}, \quad (5.110)$$

and thus

$$P_e \leq 2 + \frac{1}{M} \left(-2\text{Tr}(G) + \min\{\text{Tr}(G^2 - G), 2\text{Tr}(G)\} \right). \quad (5.111)$$

Assume that the M -codewords are independently generated according to the distribution $P_{\underline{X}}(\underline{x}) = \prod_{i=1}^N P_X(x_i)$. The expectations of G and $G^2 - G$ over the random

code are

$$\begin{aligned}
\mathbb{E}[G] &= \sum_{j=1}^M \mathbb{E}[|\psi_{f(j)}\rangle\langle\psi_{f(j)}|] = M \left(\sum_x P_X(x) |\psi_x\rangle\langle\psi_x| \right)^{\otimes N} = M\rho^{\otimes N}, \\
\mathbb{E}[G^2 - G] &= \mathbb{E} \left[\sum_{j,k=1}^M |\psi_{f(j)}\rangle\langle\psi_{f(j)}| |\psi_{f(k)}\rangle\langle\psi_{f(k)}| - \sum_{j=1}^M |\psi_{f(j)}\rangle\langle\psi_{f(j)}| \right] \\
&= \sum_{j=1}^M \sum_{k \neq j} \mathbb{E}[|\psi_{f(j)}\rangle\langle\psi_{f(j)}| |\psi_{f(k)}\rangle\langle\psi_{f(k)}|] \\
&= M(M-1)(\rho^{\otimes N})^2.
\end{aligned} \tag{5.112}$$

The expected P_e over the random code is then

$$\mathbb{E}[P_e] \leq 2 - 2\text{Tr}(\rho^{\otimes N}) + \min \left\{ (M-1)\text{Tr} \left((\rho^{\otimes N})^2 \right), 2\text{Tr}(\rho^{\otimes N}) \right\}. \tag{5.113}$$

When we denote the eigenvalues of ρ as σ_m , $m = 1, \dots, r$, where r is the rank of ρ , the eigenvalues of $\rho^{\otimes N}$, which will be denoted as $\boldsymbol{\sigma}_n$, $n = 1, \dots, r^N$, are product of N numbers each of which is chosen from $\{\sigma_1, \dots, \sigma_r\}$. There are total r^N such combinations, so that the number of eigenvalues of $\rho^{\otimes N}$ is r^N . By using the notations of the eigenvalues,

$$\mathbb{E}[P_e] \leq 2 - 2 \sum_{n=1}^{r^N} \boldsymbol{\sigma}_n + \sum_{n=1}^{r^N} \min \left\{ (M-1) (\boldsymbol{\sigma}_n)^2, 2\boldsymbol{\sigma}_n \right\}. \tag{5.114}$$

Now, by using $\min\{x, y\} \leq x^s y^{1-s}$, for $x, y \geq 0$ and every $s \in [0, 1]$,

$$\min \left\{ (M-1) (\boldsymbol{\sigma}_n)^2, 2\boldsymbol{\sigma}_n \right\} \leq (M-1)^s 2^{1-s} \boldsymbol{\sigma}_n^{1+s} \leq 2(M-1)^s \boldsymbol{\sigma}_n^{1+s}, \tag{5.115}$$

and hence

$$\mathbb{E}[P_e] \leq 2 - 2 \sum_{n=1}^{r^N} \boldsymbol{\sigma}_n + \sum_{n=1}^{r^N} 2(M-1)^s \boldsymbol{\sigma}_n^{1+s}. \tag{5.116}$$

Now from $\text{Tr}(\rho^{\otimes N}) = \sum_{n=1}^{r^N} \sigma_n = 1$,

$$\mathbb{E}[P_e] \leq 2(M-1)^s \sum_{n=1}^{r^N} \sigma_n^{1+s}. \quad (5.117)$$

Finally by using

$$\sum_{n=1}^{r^N} \sigma_n^{1+s} = \left(\sum_{m=1}^r \sigma_m^{1+s} \right)^N = (\text{Tr}(\rho^{1+s}))^N, \quad (5.118)$$

we have

$$\mathbb{E}[P_e] \leq 2(M-1)^s (\text{Tr}(\rho^{1+s}))^N \leq 2 \exp \left[-N(-\log \text{Tr}(\rho^{1+s}) - sR) \right], \quad (5.119)$$

since $M-1 \leq e^{NR}$. This bound is true for all input distributions P_X and every $s \in [0, 1]$, so that

$$\mathbb{E}[P_e] \leq 2e^{-NE(R)} \quad (5.120)$$

for $E(R)$ in (5.107). Moreover, $\mathbb{E}[P_e] \leq 2e^{-NE(R)}$ over the random code implies an existence of a code whose average error probability $P_e \leq 2e^{-NE(R)}$. It concludes the proof of Lemma 5.4.

Conclusion

■ 6.1 Summary of Main Contributions

Due to the peculiar properties of quantum mechanics, such as the no-cloning theorem and the non-reversible measurement process, the extraction of classical information from quantum states faces new challenges that are not encountered in classical information processing. To calculate the fundamental limits of communication efficiency in quantum channels, it is usually assumed in previous theories that a large number of quantum states can be collectively measured at one time. However, this assumption, in fact, becomes the primary barrier that prevents practical implementation of capacity-achieving joint detection receivers.

The purpose of this thesis is to study the performance limits of quantum channels under practical assumptions of quantum receivers, by investigating the fundamental question of how to design the measurement process to efficiently extract information from quantum states, when the possible types of measurements are restricted to particular sets of practically implementable quantum receivers.

In Chapter 3, we consider *adaptive measurements*, with which we measure each received quantum state one at a time, and then update the next measurement process based on the previous observations. We analyze the performance of adaptive measurements for quantum detection problems. We derive the necessary and sufficient conditions for adaptive measurement to perform as well as the optimal entangling measurement that achieves the theoretical lower bound of detection error probability, the

Helstrom limit. We show that for binary hypothesis testing (BHT), the greedy algorithm that minimizes the detection error probability at each instant, with the updated posterior probabilities, can meet the necessary and sufficient conditions for the optimum adaptive measurement. We show that the Dolinar receiver, which has been known to perform optimally for the BHT between two coherent states, indeed is a physical translation of the optimal adaptive measurement.

In Chapter 4, we provide one more different viewpoint to derive the Dolinar receiver. We show that for the binary hypothesis testing between two ideal laser light pulses, if we update the adaptive measurement to maximize the communication efficiency at each instant, based on recursively updated knowledge of the receiver, then we can perform as well as in the case when we can collectively measure the received laser light of an entire duration in one shot. In other words, for the BHT, the adaptive measurement that maximizes the communication efficiency in each instant also minimizes the detection error probability at that moment. Using this viewpoint, we give a natural generalization of the design to general M -ary hypothesis testing problems.

We also analyze the information capacity with adaptive measurement, and compare the result with that of direct detection receivers and of arbitrary quantum receivers (the Holevo limit), using the appropriate scaling laws in the low photon number regime. Our analysis shows that if we measure each state one at a time, we cannot approach the ultimate capacity of quantum channels, the Holevo limit, even with the capability to use the previous observations to update the measurement process.

Finally, in Chapter 5, we analyze superadditivity—the phenomenon that the maximum accessible information per channel use increases strictly as the number of channel outputs jointly measured at the receiver increases—over a pure-state classical-quantum channel. We analyze the rate vs. complexity trade-off by considering the capacity of the classical discrete memoryless superchannel induced under a concatenated coding scheme, where the quantum measurement acts exclusively on the finite length inner codewords, while allowing arbitrary outer-code complexity. We prove a general lower

bound on the maximum accessible information per channel use for a finite-length joint measurement, and express it in terms of V , the quantum version of channel dispersion, and C , the channel capacity. The superadditivity is observed even in the channel capacity of a classical discrete memoryless channel (DMC) in a concatenated coding scheme due to loss of information from hard-decisions by the inner decoder over blocklength N . Under this observation, we develop a unifying framework in which superadditivity in capacity of the classical DMC and that of a classical-quantum channel—in the above sense—can both be expressed by a parameter V/C^2 , a quantity that we show is proportional to the inner-decoder measurement length N that is sufficient to achieve a given fraction α of the capacity.

The analysis and new insights into the measurement process of quantum states that we develop in this thesis can be used to improve current quantum optical communication systems and to help understand the fundamental mechanisms of the extraction of information from quantum channels.

■ 6.2 Suggestions for Future Research

We conclude by discussing some possible extensions of the work presented in this thesis.

■ 6.2.1 Adaptive Measurements for M -ary Hypothesis Testing

In Chapter 3.3, we derive the necessary and sufficient conditions for adaptive measurements to achieve the Helstrom bound. These conditions imply that the optimum adaptive measurement should guarantee the same quality of decision in terms of probability of error, for every output sequence that belongs to the same decision set, as shown in Lemma 3.3. For binary hypothesis testing between multiple-copy states, a greedy algorithm, which minimizes the detection error probability from the view of hard-decision at the current stage, combined with posterior updating, can achieve these conditions. Moreover, the Dolinar receiver, which adds a feedback control signal to the received quantum state, and measures the combined signal with a photon counter, is

an exact physical translation of the optimum adaptive measurement.

Then, how can we generalize this result, and find the optimum adaptive measurement for M -ary hypothesis testing? For the binary case, since the posterior probabilities over two hypotheses stay in a single-dimensional probability space, it is not hard to balance the quality of decision, for two possible outputs of adaptive measurement, by adjusting the adaptive measurement based on posterior updating at each stage. Moreover, the control signal of the Dolinar receiver provides enough degree of freedom to tune the adaptive measurement to guarantee the symmetric evolvement of the posterior distribution regardless of how many photons arrive at the output of the photon counter.

However, this is no longer true for more than binary hypothesis testing. Since the posterior probabilities over M -ary hypotheses stay in $(M - 1)$ -dimensional space, there are basically infinitely many directions in which the posterior distribution can evolve at each instant. To minimize the average probability of error while guaranteeing the same quality of error for every possible output sequence belonging to the same decision set, we need to find the direction on the probability space in which the posterior distribution can make the biggest move at each instant toward one of the vertices, while balancing the progress over possible outputs of measurement. Finding such a direction is a very hard optimization problem. One of the promising ways to solve such a problem is to use the local approximations of information measures such as Kullback-Leibler divergence on the probability space, as suggested in [25]. However, it is important to note that the solution from the local approximation does not always result in a globally optimum solution.

The physical implementation of the optimum adaptive measurement also encounters a new challenge for M -ary hypothesis testing. To implement an adaptive measurement that has the ability to evolve the posterior distribution to a particular direction on the $(M - 1)$ -dimensional probability space, we need a receiver that has $(M - 1)$ -parameters to adjust at each instant. However, the Dolinar type of receiver has only one parameter, i.e., the complex amplitude of control signal, to tune over time. Therefore, we need

to resort to other optical devices, for example, squeezers, to provide enough degree of freedom for a receiver to tune its adaptive measurement for more than binary hypothesis testing.

■ 6.2.2 Quantifying the Efficiency of Measurement Process: Time-Varying Metrics

In Chapter 4.3, we discuss a generalization of the Dolinar receiver for M -ary hypothesis testing. In particular, we consider the family of Rényi entropy, which is a general class of entropy to measure the efficiency of communication, to design the control signal of the Dolinar receiver. We argue that the order of Rényi entropy to optimize the control signal of the Dolinar receiver should be time-varying, and provide intuitions to choose the right order, depending on how much time is left before the final hard-decision.

To illustrate this point, we show numerical simulation results of empirical error probability for a ternary hypothesis testing problem. For simplicity, we compare two cases where the orders are fixed to be either 1 or 100 throughout the simulation time. Using these simulation results, we confirm our intuition: when we have enough time to collect information before the final decision, at the beginning of communication, it is desirable to choose a smaller order to maximize the mutual information of the channel generated by the Dolinar receiver. On the other hand, when we need to make a final decision immediately, a larger order is preferable.

It is an interesting yet challenging research problem to find and analyze the proper time-varying measure with which we can optimize the efficiency for the measurement process of quantum states. For example, finding the optimum order of Rényi entropy at each instant is very hard, since we need to quantify how the currently collected information will affect the quality of the final decision. One potential starting point might be to select a set of discrete numbers for the possible candidates for the orders of Rényi entropy and also to divide the simulation time into a few “phases” of communication, depending on how much time is left before the final hard-decision. Then, with this

discrete setup, we can use dynamic programming to find the best sequence of orders of Rényi entropy that gives the minimum average probability of error among all the possible combinations. We can compare the resulting average probability of error with the Helstrom limit. It would also be interesting to consider the asymptotic decreasing rate of average probability of error as time goes to infinity, and compare the results under this asymptotic regime.

■ 6.2.3 Adaptive Sampling/Querying for Classical Inference Problems

The insights that we provide in Chapters 3 and 4 can also be applied for adaptive sampling and querying in classical inference problems. For example, let us consider a target localization problem with adaptive querying. Assume that we want to find the location of a target in d -dimensional space by querying an oracle who knows the exact location of the target. The type of question we can ask an oracle is whether or not the target is located within a particular region of the d -dimensional space. We observe a noisy version of the oracle's answer. For example, when the oracle's answer is 1, which indicates "yes," we observe 1 with probability p , and 0, which indicates "no," with probability $1-p$. When the oracle's answer is 0, we observe 1 with q and 0 with $1 - q$. Our goal is to minimize the mean squared error of the target location after a fixed number of queries. This kind of problem has been widely studied in previous literatures [5, 26, 44].

We want to find the best adaptive querying strategy that minimizes the mean squared error. Based on the previous noisy answers from the oracle, we need to design the next query to extract as much useful information as possible to minimize the squared error by using the oracle's next answer. Even though the final goal is to minimize the mean squared error, this cost function is hard to track and analyze in the course of adaptive querying. Consequently, in previous literatures, it has instead been widely adopted to design the adaptive querying to minimize the conditional entropy of posterior distribution over the target location. However, as pointed out in Chapter 4.3,

minimizing the conditional entropy of posterior distribution does not always give the best performance in terms of minimizing the probability of error or other general error functions.

A better way to approach this problem is thus to use the time-varying metrics in designing the query, depending on how many queries we have a chance to address to the oracle before the final decision is made. We can again use the family of Rényi entropy with different orders. When we have enough chances to ask the oracle about the target location, at the very beginning, it would be preferable to design queries with which we can maximize the mutual information between the posterior distribution of the target location and the noisy answer from the oracle. However, as the deadline for the final localization approaches, it would be preferable to choose a larger order of Rényi entropy to concentrate the measure of the posterior distribution in a more localized region.

Therefore, the insights that we develop in Chapter 4.3 and further discussions in Chapter 6.2.2 can be applied to much more general setups of not only quantum, but also classical adaptive inference problems.

■ 6.2.4 Finite Blocklength Joint Measurement: Converse Bounds

In Chapter 5.4, we provide a lower bound on the maximum achievable information rate, C_N , of a quantum channel $W \rightarrow x \rightarrow |\psi_x\rangle$, $x \in \mathcal{X}$, at a finite blocklength N of quantum measurements under concatenated coding. To calculate the exact C_N , we need to find the best superchannel that can be generated by a finite blocklength inner code-joint measurement pair. Since the complexity of this optimization problem increases exponentially with the blocklength N , instead of trying to calculate the exact C_N , we provide a lower bound on C_N in Theorem 5.2.

The proof of this lower bound is based on two ideas: First, instead of tracking the exact superchannel distribution, which depends on the detailed structure of the length- N inner code and joint measurement, we focus on one representative quantity, the average probability of error of the inner code, p_e , calculated from the superchannel

distribution, that can be easily analyzed and optimized. Second, among superchannels that have the same value of p_e , we find a superchannel whose mutual information is the smallest.

The natural question to ask next is how to find an upper bound on C_N in the finite regime of N , and how close the lower and upper bounds would be. One promising approach is again to focus on the representative quantity, p_e , of superchannels, and then to find a superchannel whose mutual information is the largest among those having the same value of p_e . In [11], the so-called “telltale” superchannel distribution is provided, whose mutual information is the largest among every superchannel (with the same p_e) that satisfies the symmetry in its channel distribution between every input, and also between every output except the right estimate. Because of this assumption of symmetry, the telltale superchannel might not be the superchannel of the largest mutual information among every superchannel with the same p_e . However, we can use the characteristics of this superchannel and attempt to extend it even for non-symmetric superchannels, in order to find an upper bound of C_N .

Bibliography

- [1] Antonio Acín, Emili Bagan, Marià Baig, Ll Masanes, and R Muñoz-Tapia. Multiple-copy two-state discrimination with individual measurements. *Physical Review A*, 71(3):032338, 2005.
- [2] Antonio Assalini, Nicola Dalla Pozza, and Gianfranco Pierobon. Revisiting the Dolinar receiver through multiple-copy state discrimination theory. *Physical Review A*, 84(2):022342, 2011.
- [3] Julio T Barreiro, Tzu-Chieh Wei, and Paul G Kwiat. Beating the channel capacity limit for linear photonic superdense coding. *Nature physics*, 4(4):282–286, 2008.
- [4] Charles H Bennett, David P DiVincenzo, Christopher A Fuchs, Tal Mor, Eric Rains, Peter W Shor, John A Smolin, and William K Wootters. Quantum nonlocality without entanglement. *Physical Review A*, 59(2):1070, 1999.
- [5] Rui Castro and Robert Nowak. Active learning and sampling. In *Foundations and Applications of Sensor Management*, pages 177–200. Springer, 2008.
- [6] Hye Won Chung, S Guha, and Lizhong Zheng. On capacity of optical channels with coherent detection. In *Information Theory Proceedings (ISIT), 2011 IEEE International Symposium on*, pages 284–288. IEEE, 2011.
- [7] Imre Csiszar and János Körner. *Information theory: coding theorems for discrete memoryless systems*. Cambridge University Press, 2011.

-
- [8] Marco Dalai. Sphere packing bound for quantum channels. In *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, pages 160–164. IEEE, 2012.
- [9] Edward Davies. Information and quantum measurement. *Information Theory, IEEE Transactions on*, 24(5):596–599, 1978.
- [10] Samuel Joseph Dolinar. An optimum receiver for the binary coherent state quantum channel. *MIT Research Laboratory of Electronics Quarterly Progress Report*, 111:115–120, 1973.
- [11] G David Forney. *Concatenated codes*, volume 11. Citeseer, 1966.
- [12] Christopher A Fuchs and Carlton M Caves. Mathematical techniques for quantum communication theory. *Open Systems & Information Dynamics*, 3(3):345–356, 1995.
- [13] Robert G Gallager. *Information theory and reliable communication*. 1968.
- [14] Vittorio Giovannetti, Saikat Guha, Seth Lloyd, Lorenzo Maccone, Jeffrey H Shapiro, and Horace P Yuen. Classical capacity of the lossy bosonic channel: The exact solution. *Physical review letters*, 92(2):027902, 2004.
- [15] Vittorio Giovannetti, Seth Lloyd, and Lorenzo Maccone. Achieving the Holevo bound via sequential measurements. *Physical Review A*, 85(1):012302, 2012.
- [16] Saikat Guha. Structured optical receivers to attain superadditive capacity and the Holevo limit. *Physical Review Letters*, 106(24):240502, 2011.
- [17] Saikat Guha and Mark M Wilde. Polar coding to achieve the Holevo capacity of a pure-loss optical channel. In *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, pages 546–550. IEEE, 2012.
- [18] Saikat Guha, Zachary Dutton, and Jeffrey H Shapiro. On quantum limit of optical communications: concatenated codes and joint-detection receivers. In *Information*

-
- Theory Proceedings (ISIT), 2011 IEEE International Symposium on*, pages 274–278. IEEE, 2011.
- [19] Matthew B Hastings. Superadditivity of communication capacity using entangled inputs. *Nature Physics*, 5(4):255–257, 2009.
- [20] Paul Hausladen, Richard Jozsa, Benjamin Schumacher, Michael Westmoreland, and William K. Wootters. Classical information capacity of a quantum channel. *Phys. Rev. A*, 54:1869–1876, Sep 1996. doi: 10.1103/PhysRevA.54.1869.
- [21] Carl W Helstrom et al. *Quantum detection and estimation theory*, volume 84. Academic press New York, 1976.
- [22] Alexander S Holevo. Bounds for the quantity of information transmitted by a quantum communication channel. *Problemy Peredachi Informatsii*, 9(3):3–11, 1973.
- [23] Alexander S Holevo. The capacity of the quantum channel with general signal states. *Information Theory, IEEE Transactions on*, 44(1):269–273, 1998. ISSN 0018-9448. doi: 10.1109/18.651037.
- [24] Alexander S Holevo. Coding theorems for quantum channels. *arXiv preprint quant-ph/9809023*, 1998.
- [25] Shao-Lun Huang and Lizhong Zheng. Linear information coupling problems. In *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, pages 1029–1033. IEEE, 2012.
- [26] Bruno Jedynek, Peter I Frazier, Raphael Sznitman, et al. Twenty questions with noise: Bayes optimal policies for entropy loss. *Journal of Applied Probability*, 49(1):114–136, 2012.
- [27] Robert S Kennedy. On the optimum receiver for the m-ary linearly independent pure state problem. *MIT Res. Lab. Electron. Quart. Prog. Rep*, 110:142–146, 1973.

-
- [28] Robert S Kennedy. Uniqueness of the optimum receiver for the m-ary pure state problem. *MIT Res. Lab. Electron. Quart. Prog. Rep.*, 113:129–130, 1974.
- [29] Amos Lapidoth, Jeffrey H Shapiro, Vinodh Venkatesan, and Ligong Wang. The discrete-time Poisson channel at low input powers. *Information Theory, IEEE Transactions on*, 57(6):3260–3272, 2011.
- [30] William Matthews and Stephanie Wehner. Finite blocklength converse bounds for quantum channels. *arXiv preprint arXiv:1210.4722*, 2012.
- [31] Michael A Nielsen and Isaac L Chuang. *Quantum computation and quantum information*. Cambridge university press, 2010.
- [32] Michael Nussbaum and Arleta Szkoła. The chernoff lower bound for symmetric quantum hypothesis testing. *The Annals of Statistics*, pages 1040–1057, 2009.
- [33] Asher Peres and William K Wootters. Optimal detection of quantum information. *Physical Review Letters*, 66(9):1119–1122, 1991.
- [34] Yury Polyanskiy, H Vincent Poor, and Sergio Verdú. Channel coding rate in the finite blocklength regime. *Information Theory, IEEE Transactions on*, 56(5):2307–2359, 2010.
- [35] Masahide Sasaki, Kentaro Kato, Masayuki Izutsu, and Osamu Hirota. A demonstration of superadditivity in the classical capacity of a quantum channel. *Physics Letters A*, 236(1):1–4, 1997.
- [36] Masahide Sasaki, Kentaro Kato, Masayuki Izutsu, and Osamu Hirota. Quantum channels showing superadditivity in classical capacity. *Physical Review A*, 58(1):146, 1998.
- [37] Benjamin Schumacher and Michael D Westmoreland. Sending classical information via noisy quantum channels. *Physical Review A*, 56:131–138, 1997.

-
- [38] Shlomo Shamai. Capacity of a pulse amplitude modulated direct detection photon channel. In *Communications, Speech and Vision, IEE Proceedings I*, volume 137, pages 424–430. IET, 1990.
- [39] Claude E Shannon, Robert G Gallager, and Elwyn R Berlekamp. Lower bounds to error probability for coding on discrete memoryless channels. i. *Information and Control*, 10(1):65–103, 1967.
- [40] Peter W Shor. The adaptive classical capacity of a quantum channel, or information capacities of three symmetric pure states in three dimensions. *IBM Journal of Research and Development*, 48(1):115–137, 2004.
- [41] Masaki Sohma and Osamu Hirota. Binary discretization for quantum continuous channels. *Physical Review A*, 62(5):52312, 2000.
- [42] Masahiro Takeoka, Hari Krovi, and Saikat Guha. Achieving the Holevo capacity of a pure state classical-quantum channel via unambiguous state discrimination. In *Information Theory Proceedings (ISIT), 2013 IEEE International Symposium on*, pages 166–170. IEEE, 2013.
- [43] Marco Tomamichel and Vincent YF Tan. Second-order asymptotics of classical-quantum channels. *arXiv preprint arXiv:1308.6503*, 2013.
- [44] Theodoros Tsiligkaridis, Brian M Sadler, and Alfred O Hero. A collaborative 20 questions model for target search with human-machine interaction. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 6516–6520. IEEE, 2013.
- [45] Sergio Verdú. Spectral efficiency in the wideband regime. *Information Theory, IEEE Transactions on*, 48(6):1319–1343, 2002.
- [46] Mark M Wilde and Saikat Guha. Polar codes for classical-quantum channels. 2011.

-
- [47] Mark M Wilde, Saikat Guha, Si-Hui Tan, and Seth Lloyd. Explicit capacity-achieving receivers for optical communication and quantum reading. In *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, pages 551–555. IEEE, 2012.
- [48] Mark M Wilde, Olivier Landon-Cardinal, and Patrick Hayden. Towards efficient decoding of classical-quantum polar codes. *arXiv preprint arXiv:1302.0398*, 2013.
- [49] Aaron D Wyner. Capacity and error exponent for the direct detection photon channel. ii. *Information Theory, IEEE Transactions on*, 34(6):1462–1471, 1988.
- [50] Horace P Yuen and Masanao Ozawa. Ultimate information carrying limit of quantum systems. *Physical Review Letters*, 70:363–366, 1993.
- [51] Horace P Yuen, Robert S Kennedy, and Melvin Lax. Optimum testing of multiple hypotheses in quantum detection theory. *Information Theory, IEEE Transactions on*, 21(2):125–134, 1975.
- [52] Lizhong Zheng, David NC Tse, and Muriel Médard. Channel coherence in the low-snr regime. *Information Theory, IEEE Transactions on*, 53(3):976–997, 2007.