

# Authentication of Quantum Messages

Howard Barnum<sup>\*</sup>, Claude Crépeau<sup>†</sup>, Daniel Gottesman<sup>‡</sup>, Adam Smith<sup>§</sup>, and Alain Tapp<sup>¶</sup>

## Abstract

Authentication is a well-studied area of classical cryptography: a sender  $\mathcal{A}$  and a receiver  $\mathcal{B}$  sharing a classical private key want to exchange a classical message with the guarantee that the message has not been modified or replaced by a dishonest party with control of the communication line. In this paper we study the authentication of messages composed of *quantum states*.

We give a formal definition of authentication in the quantum setting. Assuming  $\mathcal{A}$  and  $\mathcal{B}$  have access to an insecure quantum channel and share a private, classical random key, we provide a *non-interactive* scheme that both enables  $\mathcal{A}$  to encrypt and authenticate (with unconditional security) an  $m$  qubit message by encoding it into  $m + s$  qubits, where the probability decreases exponentially in the security parameter  $s$ . The scheme requires a private key of size  $2m + O(s)$ . To achieve this, we give a highly efficient protocol for testing the purity of shared EPR pairs.

It has long been known that learning information about a general quantum state will necessarily disturb it. We refine this result to show that such a disturbance can be done with few side effects, allowing it to circumvent cryptographic protections. Consequently, any scheme to authenticate quantum messages must also encrypt them. In contrast, no such constraint exists classically: authentication and encryption are independent tasks, and one can authenticate a message while leaving it publicly readable.

This reasoning has two important consequences: On one hand, it allows us to give a lower bound of  $2m$  key bits for authenticating  $m$  qubits, which makes our protocol asymptotically optimal. On the other hand, we use it to show that digitally signing quantum states is impossible, even with only computational security.

**Keywords.** Authentication, quantum information.

---

<sup>\*</sup> CCS-3 Group, Los Alamos National Laboratories, Los Alamos, New Mexico 87554 USA. e-mail: barnum@lanl.gov. Supported by the US DOE; part of this research was done while working at the University of Bristol (UK), and supported by the EU QAIP Consortium (IST-1999-11234).

<sup>†</sup> School of Computer Science, McGill University, Montréal (Québec), Canada. e-mail: crepeau@cs.mcgill.ca. Supported in part by Québec's FCAR and Canada's NSERC.

<sup>‡</sup> UC Berkeley, EECS: Computer Science Division, Soda Hall 585, Berkeley, California 94720, USA. e-mail: gottesma@eecs.berkeley.edu. Supported by the Clay Mathematics Institute.

<sup>§</sup> M.I.T., Laboratory for Computer Science, 200 Technology Square, Cambridge MA 02139, USA. e-mail: asmith@theory.lcs.mit.edu. Supported in part by U.S. Army Research Office Grant DAAD19-00-1-0177. Some of this research was done while the author was visiting McGill University.

<sup>¶</sup> Département IRO, Université de Montréal, C.P. 6128, succursale centre-ville, Montréal (Québec), Canada H3C 3J7. e-mail: tappa@iro.umontreal.ca. Part of this research was done while working at Department of Combinatorics and Optimization, University of Waterloo and McGill University.

# 1 Introduction

Until recently, the expression “quantum cryptography” referred mostly to quantum key distribution protocols [5, 4, 13]. However, these words now refer to a larger set of problems. While QKD and many other quantum protocols attempt to provide improved security for tasks involving classical information, an emerging area of quantum cryptography attempts instead to create secure protocols for tasks involving *quantum* information. One standard cryptographic task is the *authentication* of messages:  $\mathcal{A}$  transmits some information to  $\mathcal{B}$  over an insecure channel, and they wish to be sure that it has not been tampered with *en route*. When the message is classical, and  $\mathcal{A}$  and  $\mathcal{B}$  share a random private key, this problem can be solved by, for instance, the Wegman-Carter scheme [12]. In this paper, we discuss the analogous question for quantum messages.

**A naive approach** If we assume  $\mathcal{A}$  and  $\mathcal{B}$  share a private *quantum* key in the form of  $m$  EPR pairs, as well as some private classical key, there is a straightforward solution to this problem:  $\mathcal{A}$  simply uses quantum teleportation [6] to send her message to Bob, authenticating the  $2m$  classical bits transmitted in the teleportation protocol. If  $\mathcal{A}$  and  $\mathcal{B}$  initially share only a classical key, however, the task is more difficult. We start with a simple approach: first distribute EPR pairs (which might get corrupted in transit), and then use entanglement purification [7] to establish clean pairs for teleportation. This can be improved: we do not need a full-scale entanglement purification protocol, which produces good EPR pairs even if the channel is noisy; instead we only need something we call a *purity testing protocol*, which checks that EPR pairs are correct, but does not attempt to repair them in case of error.

Unfortunately, any such protocol will have to be interactive, since  $\mathcal{A}$  must first send some qubits to  $\mathcal{B}$  and then wait for confirmation of receipt before completing the transmission. This is unsuitable for situations in which a message is stored and must be checked for authenticity at a later time. Also, this interactive protocol achieves something stronger than what is required of a quantum authentication scheme: at the end of the purity-testing based scheme, *both* Alice and Bob know that the transmission was successful, whereas for authentication, we only require that Bob knows.

**Contributions** In this paper we study *non-interactive* quantum authentication schemes with classical keys. Our primary contributions are:

- Formal definition of authentication for quantum states

In classical authentication, one simply limits the probability that the adversary can make *any* change to the state without detection. This condition is too stringent for quantum information, where we only require high fidelity to the original state. We state our definition in terms of the transmission of pure states (section 3), but also show that the same definition implies security for mixed or entangled states.

- Construction of efficient purity testing protocols

We show how to create purity-testing protocols using families of quantum error-correcting codes with a particular covering property, namely that any Pauli error is detected by most of the codes in the family. We construct an efficient such family based on projective geometry, yielding a purity-testing protocol requiring only  $O(s)$  (classical) bits of communication, where  $s$  is the security parameter (section 4).

Purity-testing codes have not explicitly appeared before in the literature, but have been present implicitly in earlier work, for instance [15, 19]. To prove our purity-testing protocols secure, we use a “quantum-to-classical” reduction, due to Lo and Chau [15]. Subsequently to our work, Ambainis, Smith, and Yang [3] used our construction of purity-testing protocols in a study of more general entanglement extraction procedures.

- Construction of non-interactive quantum authentication schemes (QAS)

We show that a secure non-interactive QAS can be constructed from any purity-testing protocol derived, as above, from QECCs (section 5). In particular, for our family of codes, we obtain an authentication scheme which requires sending  $m + s$  qubits, and consuming  $2m + O(s)$  bits of classical key for a message of  $m$  qubits. The proof techniques in the Shor and Preskill paper [19] serve as inspiration for the transformation from an interactive purity-testing protocol to a non-interactive QAS.

- Study of the relation between encryption and authentication

One feature of our authentication protocol is that it completely encrypts the quantum message being sent. We show that this is a necessary feature of *any* QAS (section 6), in striking contrast to the situation for classical information, where common authentication schemes leave the message completely intelligible. It therefore follows that any authentication protocol for an  $m$ -qubit message must use nearly  $2m$  bits of classical key, enough to encrypt the message. The protocol we present approaches this bound asymptotically.

- Impossibility of digitally signing quantum states

Since authentication requires encryption, it is impossible to create digital signature schemes for quantum messages: any protocol which allows one recipient to read a message also allows him or her to modify it without risk of detection, and therefore all potential recipients of an authenticated message must be trustworthy (section 7). This conclusion holds true even if we require only computationally secure digital signatures. Note that this does not in any way preclude the possibility of signing *classical* messages with or without quantum states [14].

Why should we prefer a scheme with classical keys to a scheme with entangled quantum keys? The task of authenticating quantum data is only useful in a scenario where quantum information can be reliably stored, manipulated, and transmitted over communication lines, so it would not be unreasonable to assume quantum keys. However, many manipulations are easier with classical keys. Certainly, the technology for storing and manipulating them is already available, but there are additional advantages. Consider, for example, public key cryptography; it is possible to sign and encrypt classical key bits with public key systems, but signing a general quantum state is impossible. Thus, quantum keys would be unsuitable for an asymmetric quantum authentication scheme such as the one we describe in section 5.1.

## 2 Preliminaries

### 2.1 Classical Authentication

In the classical setting, an authentication scheme is defined by a pair of functions  $A : \mathcal{K} \times M \rightarrow C$  and  $B : \mathcal{K} \times C \rightarrow M \times \{\text{valid}, \text{invalid}\}$  such that for any message  $\mu \in M$  and key  $k \in \mathcal{K}$  we have *completeness*

$$B_k(A_k(\mu)) = \langle \mu, \text{valid} \rangle$$

and that for any opponent algorithm  $O$ , we have *soundness*

$$\text{Prob} \{B_k(O(A_k(\mu))) \in \{\langle \mu, \text{valid} \rangle\} \cup \{\langle \mu', \text{invalid} \rangle | \mu' \in M\}\} \geq 1 - 2^{-\Omega(t)}$$

where  $t = \lg \#C - \lg \#M$  is the security parameter creating the tradeoff between the expansion of the messages and the security level. Note that we only consider information-theoretically secure schemes, not schemes that are based on computational assumptions.

Wegman and Carter [12] introduced several constructions for such schemes; their most efficient uses keys of size only  $4(t + \lg \lg m) \lg m$  and achieves security  $1 - 2^{-t+2}$ . This compares rather well to the known lower bound of  $t + \lg m - \lg t$  for such a result [12]. The same work also introduced a technique to re-use an authentication function several times by using one-time-pad encryption on the tag, so that an opponent cannot learn *anything* about the particular key being used by  $\mathcal{A}$  and  $\mathcal{B}$ . Thus, at a marginal cost of only  $t$  secret key bits per authentication, the confidentiality of the authentication function  $h$  is guaranteed and thus may be re-used (a polynomial number of times).

For the remainder of this paper, we assume the reader is familiar with the basic notions and notation of quantum computing. These can be found in textbooks such as [16]. Since we rely heavily on terminology and techniques from quantum error correction (especially stabilizer codes), appendix A provides a summary of the relevant notions.

## 2.2 Purification and purity testing

Quantum error-correcting codes (QECCs) may be used for *entanglement purification* ([7]). In this setting,  $\mathcal{A}$  and  $\mathcal{B}$  share some Bell states (say  $|\Phi^+\rangle = |00\rangle + |11\rangle$ ) which have been corrupted by transmission through a noisy quantum channel. They want a protocol which processes these imperfect EPR pairs and produces a smaller number of higher-quality pairs. We assume that  $\mathcal{A}$  and  $\mathcal{B}$  have access to an authenticated, public classical channel. At the end of the protocol, they either accept or reject based on any inconsistencies they have observed. As long as  $\mathcal{A}$  and  $\mathcal{B}$  have a noticeable probability of accepting, then conditioned on accepting, the state they share should have fidelity almost 1 to the pure state  $|\Phi^+\rangle^{\otimes m}$ . Moreover, small amounts of noise in their initial shared state should not cause failure of the protocol.

Stabilizer codes can be particularly useful for purification because of the following observation: for any stabilizer code  $Q$ , if we measure the syndrome of one half of a set of Bell states  $|\Phi^+\rangle^{\otimes n}$  and obtain the result  $y$ , then the result is the state  $|\Phi^+\rangle^{\otimes m}$ , with each of its two halves encoded in the coset with syndrome  $y$ . (Moreover, in this case the distribution on  $y$  is uniform.) If the original state is erroneous,  $\mathcal{A}$  and  $\mathcal{B}$  will likely find different syndromes, which will differ by the syndrome associated with the actual error.

Most purification protocols based on stabilizer codes require efficient error correction; we measure the syndrome, and use that information to efficiently restore the encoded state. However, one can imagine a weaker task in which Alice and Bob only want to *test* their EPR pairs for purity, i.e. they want a guarantee that if their pairs pass the test, their shared state will probably be close to  $|\Phi^+\rangle^{\otimes m}$ . In that case, we can use the code for error detection, not correction, and need only be able to encode and decode efficiently from the space  $Q$ .

## 2.3 Encryption of Quantum Messages

A useful ingredient for much recent work in quantum cryptography is the concept of quantum teleportation, put forward by Bennett et al [6]. After  $\mathcal{A}$  and  $\mathcal{B}$  have shared a singlet state,  $\mathcal{A}$  can later secretly send a single qubit in an arbitrary quantum state  $\rho$  to  $\mathcal{B}$  by measuring her half of the singlet state together with her state  $\rho$  in the Bell basis to get two classical bits  $b_0, b_1$ . As a result,  $\mathcal{B}$ 's half of the singlet state will become one of four possibilities  $\rho' := \sigma_z^{b_0} \sigma_x^{b_1} \rho \sigma_x^{b_1} \sigma_z^{b_0}$ . If  $\mathcal{A}$  sends  $b_0, b_1$ , then  $\mathcal{B}$  can easily recover  $\rho$ .

Now without the bits  $b_0, b_1$ , the state  $\rho'$  reveals no information about  $\rho$ . Thus, one can turn this into an encryption scheme which uses only a classical key: after  $\mathcal{A}$  and  $\mathcal{B}$  have secretly shared two classical bits  $b_0, b_1$ ,  $\mathcal{A}$  can later secretly send a single qubit in an arbitrary quantum state  $\rho$  to  $\mathcal{B}$  by sending him a qubit in state  $\rho'$  as above. This is called a quantum one-time pad (QOTP). This scheme is optimal [1, 9]: any quantum encryption (with a classical key) must use 2 bits of key for every transmitted qubit.

### 3 Quantum Authentication

At an intuitive level, a quantum authentication scheme is a keyed system which allows  $\mathcal{A}$  to send a state  $\rho$  to  $\mathcal{B}$  with a guarantee: if  $\mathcal{B}$  accepts the received state as “good”, the fidelity of that state to  $\rho$  is almost 1. Moreover, if the adversary makes no changes,  $\mathcal{B}$  should always accept, and the fidelity should be exactly 1.

Of course, this informal definition is impossible to attain. The adversary might always replace  $\mathcal{A}$ ’s transmitted message with a completely mixed state. There would nonetheless be a small probability that  $\mathcal{B}$  would accept, but even when he did accept, the fidelity of the received state to  $\mathcal{A}$ ’s initial state would be very low.

The problem here is that we are conditioning on  $\mathcal{B}$ ’s acceptance of the received state; this causes trouble if the adversary’s a priori chances of cheating are high. A more reasonable definition would require a tradeoff between  $\mathcal{B}$ ’s chances of accepting, and the expected fidelity of the received system to  $\mathcal{A}$ ’s initial state given his acceptance: as  $\mathcal{B}$ ’s chance of accepting increases, so should the expected fidelity.

It turns out that there is no reason to use both the language of probability and that of fidelity here: for classical tests, fidelity and probability of acceptance coincide. With this in mind we first define what constitutes a quantum authentication scheme, and then give a definition of security:

**Definition 1** A quantum authentication scheme (QAS) is a pair of polynomial time quantum algorithms  $A$  and  $B$  together with a set of classical keys  $\mathcal{K}$  such that:

- $A$  takes as input an  $m$ -qubit message system  $M$  and a key  $k \in \mathcal{K}$  and outputs a transmitted system  $T$  of  $m + t$  qubits.
- $B$  takes as input the (possibly altered) transmitted system  $T'$  and a classical key  $k \in \mathcal{K}$  and outputs two systems: a  $m$ -qubit message state  $M$ , and a single qubit  $V$  which indicates acceptance or rejection. The classical basis states of  $V$  are called  $|\text{ACC}\rangle, |\text{REJ}\rangle$  by convention.

For any fixed key  $k$ , we denote the corresponding super-operators by  $A_k$  and  $B_k$ .

Note that  $\mathcal{B}$  may well have measured the qubit  $V$  to see whether or not the transmission was accepted or rejected. Nonetheless, we think of  $V$  as a qubit rather than a classical bit since it will allow us to describe the joint state of the two systems  $M, V$  with a density matrix.

There are two conditions which should be met by a quantum authentication protocol. On the one hand, in the absence of intervention, the received state should be the same as the initial state and  $\mathcal{B}$  should accept.

On the other hand, we want that when the adversary does intervene,  $\mathcal{B}$ ’s output systems have high fidelity to the statement “either  $\mathcal{B}$  rejects or his received state is the same as that sent by  $\mathcal{A}$ ”. One difficulty with this is that it is not clear what is meant by “the same state” when  $\mathcal{A}$ ’s input is a mixed state. It turns out that it is sufficient to define security in terms of pure states; one can deduce an appropriate statement about fidelity of mixed states (see Appendix B).

Given a pure state  $|\psi\rangle \in \mathcal{H}_M$ , consider the following test on the joint system  $M, V$ : output a 1 if the first  $m$  qubits are in state  $|\psi\rangle$  or if the last qubit is in state  $|\text{REJ}\rangle$  (otherwise, output a 0). The projectors corresponding to this measurement are

$$\begin{aligned} P_1^{|\psi\rangle} &= |\psi\rangle\langle\psi| \otimes I_V + I_M \otimes |\text{REJ}\rangle\langle\text{REJ}| - |\psi\rangle\langle\psi| \otimes |\text{REJ}\rangle\langle\text{REJ}| \\ P_0^{|\psi\rangle} &= (I_M - |\psi\rangle\langle\psi|) \otimes (|\text{ACC}\rangle\langle\text{ACC}|) \end{aligned}$$

We want that for all possible input states  $|\psi\rangle$  and for all possible interventions by the adversary, the expected fidelity of  $\mathcal{B}$ ’s output to the space defined by  $P_1^{|\psi\rangle}$  is high. This is captured in the following definition of security.

**Definition 2** A QAS is secure with error  $\epsilon$  for a state  $|\psi\rangle$  if it satisfies:

Completeness: For all keys  $k \in \mathcal{K}$ :  $B_k(A_k(|\psi\rangle\langle\psi|)) = |\psi\rangle\langle\psi| \otimes |\text{ACC}\rangle\langle\text{ACC}|$

Soundness: For all super-operators  $\mathcal{O}$ , let  $\rho_{Bob}$  be the state output by  $\mathcal{B}$  when the adversary's intervention<sup>1</sup> is characterized by  $\mathcal{O}$ , that is:

$$\rho_{Bob} = \mathbb{E}_k \left[ B_k(\mathcal{O}(A_k(|\psi\rangle\langle\psi|))) \right] = \frac{1}{|\mathcal{K}|} \sum_k B_k(\mathcal{O}(A_k(|\psi\rangle\langle\psi|)))$$

where “ $\mathbb{E}_k$ ” means the expectation when  $k$  is chosen uniformly at random from  $\mathcal{K}$ . The QAS has soundness error  $\epsilon$  for  $|\psi\rangle$  if:

$$\text{Tr} \left( P_1^{|\psi\rangle} \rho_{Bob} \right) \geq 1 - \epsilon$$

A QAS is secure with error  $\epsilon$  if it is secure with error  $\epsilon$  for all states  $|\psi\rangle$ .

Note that our definition of completeness assumes that the channel connecting  $\mathcal{A}$  to  $\mathcal{B}$  is noiseless in the absence of the adversary's intervention. This is in fact not a significant problem, as we can simulate a noiseless channel using standard quantum error correction.

**Interactive protocols** In the previous section, we dealt only with non-interactive quantum authentication schemes, since that is both the most natural notion, and the one we achieve in this paper. However, there is no reason to rule out interactive protocols in which  $\mathcal{A}$  and  $\mathcal{B}$  at the end believe they have reliably exchanged a quantum message. The definitions of completeness and soundness extend naturally to this setting: as before,  $\mathcal{B}$ 's final output is a pair of systems  $M, V$ , where the state space of  $V$  is spanned by  $|\text{ACC}\rangle, |\text{REJ}\rangle$ . In that case  $\rho_{Bob}$  is  $\mathcal{B}$ 's density matrix at the end of the protocol, averaged over all possible choices of shared private key and executions of the protocol. The soundness error is  $\epsilon$ , where  $\text{Tr} \left( P_1^{|\psi\rangle} \rho_{Bob} \right) \geq 1 - \epsilon$ .

## 4 Purity Testing Codes

An important tool in our proof is the notion of a *purity testing code*, which is a way for  $\mathcal{A}$  and  $\mathcal{B}$  to ensure that they share (almost) perfect EPR pairs. We shall concentrate on purity testing codes based on stabilizer QECCs.

**Definition 3** A stabilizer purity testing code with error  $\epsilon$  is a set of stabilizer codes  $\{Q_k\}$ , for  $k \in \mathcal{K}$ , such that  $\forall E_x \in E$  with  $x \neq 0$ ,  $\#\{k|x \in Q_k^\perp - Q_k\} \leq \epsilon(\#\mathcal{K})$ .

That is, for any error  $x$  in the error group, if  $k$  is chosen later at random, the probability that the code  $Q_k$  detects  $x$  is at least  $1 - \epsilon$ .

**Definition 4** A purity testing protocol with error  $\epsilon$  is a superoperator  $\mathcal{T}$  which can be implemented with local operations and classical communication, and which maps  $2n$  qubits (half held by  $\mathcal{A}$  and half held by  $\mathcal{B}$ ) to  $2m + 1$  qubits and satisfies the following two conditions:

Completeness:  $\mathcal{T}(|\Phi^+\rangle^{\otimes n}) = |\Phi^+\rangle^{\otimes m} \otimes |\text{ACC}\rangle$

<sup>1</sup>We make no assumptions on the running time of the adversary.

Soundness: Let  $P$  be the projection on the subspace spanned by  $|\Phi^+\rangle^{\otimes m} \otimes |\text{ACC}\rangle$  and  $|\psi\rangle \otimes |\text{REJ}\rangle$  for all  $|\psi\rangle$ . Then  $\mathcal{T}$  satisfies the soundness condition if for all  $\rho$ ,

$$\text{Tr}(PT(\rho)) \geq 1 - \epsilon.$$

The obvious way of constructing a purity testing protocol  $\mathcal{T}$  is to start with a purity testing code  $\{Q_k\}$ . When Alice and Bob are given the state  $\rho$ , Alice chooses a random  $k \in \mathcal{K}$  and tells it to Bob. They both measure the syndrome of  $Q_k$  and compare. If the syndromes are the same, they accept and perform the decoding procedure for  $Q_k$ ; otherwise they reject.

**Proposition 1** *If the purity testing code  $\{Q_k\}$  has error  $\epsilon$ , then  $\mathcal{T}$  is a purity testing protocol with error  $\epsilon$ .*

The proof appears in Appendix C.

## 4.1 An Efficient Purity Testing Code

Now we will give an example of a particularly efficient purity testing code. We will use the stabilizer techniques of section A, restricting to the case  $n = rs$ . We will construct a set of codes  $Q_k$  each encoding  $m = (r-1)s$  qubits in  $n$  qubits, and show that the  $Q_k$  form a purity testing code. (Note that the construction below works just as well if instead of qubits, we use registers with dimension equal to any prime power; see appendix D for details.) Using qubits in groups of  $s$  allows us to view our field  $GF(2^{2rs})$  as both a  $2r$ -dimensional vector space over  $GF(2^s)$  and a  $2rs$ -dimensional binary vector space. We need a symplectic form that is compatible with this decomposition. One possibility is

$$B(x, y) := \text{Tr}(xy^{2^{rs}}), \quad (1)$$

where  $\text{Tr}(z) = \sum_{i=0}^{2^{rs}-1} z^{2^i}$  is the standard trace function, which maps  $GF(2^{2rs})$  onto  $GF(2)$ .

We consider a *normal rational curve* in  $PG(2r-1, 2^s)$  (the projective geometry whose points are the 1-d subspaces of the  $2r$ -dimensional vector space over  $GF(2^s)$ ). (See, e.g., the excellent introductory text [8].) Such a curve is given by:

$$\Upsilon = \{[1 : y : y^2 : \dots : y^{2^r-1}], [0 : 0 : 0 : \dots : 1]\}_{y \in GF(2^s)}. \quad (2)$$

(The colon is used to separate the coordinates of a projective point, indicating that only their ratio matters.) Thus, there are  $2^s + 1$  points on the normal rational curve.

Since each ‘‘point’’ of this curve is actually a one-dimensional subspace over  $GF(2^s)$ , it can also be considered as an  $s$ -dimensional binary subspace  $Q_k$  in a vector space of dimension  $2rs = 2n$ . We will show that  $Q_k$  is totally isotropic with respect to the symplectic inner product (1), and encodes  $m = n - s$  qubits in  $n$  qubits.

**Theorem 2** *The set of codes  $Q_k$  form a stabilizer purity testing code with error*

$$\epsilon = \frac{2r}{2^s + 1}. \quad (3)$$

*Each code  $Q_k$  encodes  $m = (r-1)s$  qubits in  $n = rs$  qubits.*

Proof of this is in Appendix D.

## 5 Protocols

In this section we describe a secure non-interactive quantum authentication scheme (Protocol 5.2) which satisfies the definition of section 3.

In order to prove our scheme secure, we begin with a purity testing protocol as per Section 4 (summarized as Protocol 5.1). The security of this protocol follows from Prop. 1. We then perform several transformations to the protocol that strictly preserve its security and goals but which remove the interaction, replacing it with a shared private key. We thus obtain two less interactive intermediate protocols (Protocols E.1 and E.2) and a final protocol (Protocol 5.2), which is completely non-interactive. The transformations are similar in flavor to those of Shor and Preskill [19], who use the technique to obtain a simple proof of the security of a completely different task, namely the BB84 [5] quantum key exchange scheme.

### Protocol 5.1 ( Purity Testing Based Protocol )

- 1:  $\mathcal{A}$  and  $\mathcal{B}$  agree on some stabilizer purity testing code  $\{Q_k\}$
- 2:  $\mathcal{A}$  generates  $2n$  qubits in state  $|\Phi^+\rangle^{\otimes n}$ .  $\mathcal{A}$  sends the first half of each  $|\Phi^+\rangle$  state to  $\mathcal{B}$ .
- 3:  $\mathcal{B}$  announces that he has received the  $n$  qubits.
- 4:  $\mathcal{A}$  picks a random  $k \in \mathcal{K}$ , and announces it to  $\mathcal{B}$ .
- 5:  $\mathcal{A}$  and  $\mathcal{B}$  measure the syndrome of the stabilizer code  $Q_k$ .  $\mathcal{A}$  announces her results to  $\mathcal{B}$  who compares them to his own results. If any error is detected,  $\mathcal{B}$  aborts.
- 6:  $\mathcal{A}$  and  $\mathcal{B}$  decode their  $n$ -qubit words according to  $Q_k$ . Each is left with  $m$  qubits, which together should be nearly in state  $|\Phi^+\rangle^{\otimes m}$ .
- 7:  $\mathcal{A}$  uses her half of  $|\Phi^+\rangle^{\otimes m}$  to teleport an arbitrary  $m$ -qubit state  $\rho$  to  $\mathcal{B}$ .

Following the notation of Section 4, let  $P$  be the projector onto the subspace described by “either  $\mathcal{B}$  has aborted or the joint state held by  $\mathcal{A}$  and  $\mathcal{B}$  is  $|\Phi^+\rangle^{\otimes m}$ ”. Let  $\rho_{AB}$  be the joint density matrix of  $\mathcal{A}$  and  $\mathcal{B}$ ’s systems. Then Prop. 1 states that at the end of step 6,  $\text{Tr}(P\rho_{AB})$  is exponentially close to 1 in  $n$ . The soundness of our first authentication protocol follows immediately:

**Corollary 3** *If  $\mathcal{A}$  and  $\mathcal{B}$  are connected by an authenticated classical channel, then Protocol 5.1 is a secure interactive quantum authentication protocol, with soundness error exponentially small in  $n$ .*

The proof is straightforward; we give it explicitly in Appendix E.

**Theorem 4** *When the purity testing code  $\{Q_k\}$  has error  $\epsilon$ , the protocol 5.2 is a secure quantum authentication scheme with key length  $O(n + \log_2(\#\mathcal{K}))$  and soundness error  $\epsilon$ . In particular, for the purity testing code described in Section 4.1, the authentication scheme has key length  $2m + s + \log_2(2^s + 1) \leq 2n + 1$  and soundness error  $2n/[s(2^s + 1)]$ , where  $m$  is the length of the message in qubits,  $s$  is the security parameter, and  $\mathcal{A}$  sends a total of  $n = m + s$  qubits.*

*Proof:* From Corollary 3 we have that Protocol 5.1 is a secure interactive authentication protocol. We show that Protocol 5.2 is equivalent to Protocol 5.1, in the sense that any attack on Protocol 5.2 implies an equally successful attack on Protocol 5.1. To do so, we proceed by a series of reductions; the details appear in Appendix E.



### Protocol 5.2 ( Non-interactive authentication )

- 1:** *Preprocessing:*  $\mathcal{A}$  and  $\mathcal{B}$  agree on some stabilizer purity testing code  $\{Q_k\}$  and some private and random binary strings  $k$ ,  $x$ , and  $y$ .
- 2:**  $\mathcal{A}$   $q$ -encrypts  $\rho$  as  $\tau$  using key  $x$ .  $\mathcal{A}$  encodes  $\tau$  according to  $Q_k$  for the code  $Q_k$  with syndrome  $y$  to produce  $\sigma$ .  $\mathcal{A}$  sends the result to  $\mathcal{B}$ .
- 3:**  $\mathcal{B}$  receives the  $n$  qubits. Denote the received state by  $\sigma'$ .  $\mathcal{B}$  measures the syndrome  $y'$  of the code  $Q_k$  on his qubits.  $\mathcal{B}$  compares  $y$  to  $y'$ , and aborts if any error is detected.  $\mathcal{B}$  decodes his  $n$ -qubit word according to  $Q_k$ , obtaining  $\tau'$ .  $\mathcal{B}$   $q$ -decrypts  $\tau'$  using  $x$  and obtains  $\rho'$ .

## 5.1 Public Key Quantum Authentication

Unlike its classical counterpart, quantum information can be authenticated in a public key setting but not in a way that can be demonstrated to a judge. In section 6, we show the impossibility of a digital signature scheme for quantum information; here, we instead introduce the notion of public key quantum authentication.

Let  $E_b, D_b$  be  $\mathcal{B}$ 's public and private keyed algorithms to a PKC resistant to quantum computers' attacks. Let  $S_a, V_a$  be  $\mathcal{A}$ 's private and public keyed algorithms to a digital signature scheme resistant to quantum computers' attacks. These may be either be protocols which are secure with respect to a computational assumption [17] or with unconditional security [14]. To perform authentication,  $\mathcal{A}$  picks secret and random binary strings  $k$ ,  $x$ , and  $y$ , and uses them as keys to  $q$ -authenticate  $\rho$  as  $\rho'$ .  $\mathcal{A}$  encrypts and signs the key as  $\sigma := S_a(E_b(k|x|y))$ .  $\mathcal{A}$  sends  $(\rho', \sigma)$  to  $\mathcal{B}$ . To verify a state,  $\mathcal{B}$  verifies  $\mathcal{A}$ 's signature on  $\sigma$  using  $V_a$  and then discovers the key  $k$ ,  $x$  and  $y$  using his private decryption function  $D_b$ .  $\mathcal{B}$  checks that  $\rho'$  is a valid  $q$ -authenticated message according to key  $k$ ,  $x$ ,  $y$ , and recovers  $\rho$ .

## 6 Good Authentication Implies Good Encryption

One notable feature of any protocol derived using Theorem 4 is that the information being authenticated is also completely encrypted. For classical information, authentication and encryption can be considered completely separately, but in this section we will show that quantum information is different. While quantum states can be encrypted without any form of authentication, the converse is not true: any scheme which guarantees authenticity must also encrypt the quantum state almost perfectly.

To show this, let us consider any fixed authentication scheme. Denote by  $\rho_{|\psi\rangle}$  the density matrix transmitted in this scheme when Alice's input is  $|\psi\rangle$ . Let  $\rho_{|\psi\rangle}^{(k)}$  denote the density matrix for key  $k$ .

**Definition 5** *An encryption scheme with error  $\epsilon$  for quantum states hides information so that if  $\rho_0$  and  $\rho_1$  are any two distinct encrypted states, then the trace distance  $D(\rho_0, \rho_1) = \frac{1}{2}\text{Tr} |\rho_0 - \rho_1| \leq \epsilon$ .*

We claim that any good QAS must necessarily also be a good encryption scheme. That is:

**Theorem 5 (Main Lower Bound)** *A QAS with error  $\epsilon$  is an encryption scheme with error at most  $4\epsilon^{1/6}$ .*

**Corollary 6** *A QAS with error  $\epsilon$  requires at least  $2m(1 - \text{poly}(\epsilon))$  classical key bits.*

We prove this corollary in Appendix F. For now, we concentrate on the Theorem 5.

The intuition behind the proof of this main theorem is that measurement disturbs quantum states, so if the adversary can learn information about the state, she can change the state. More precisely, if the adversary can distinguish between two states  $|0\rangle$  and  $|1\rangle$ , she can change the state  $|0\rangle + |1\rangle$  to  $|0\rangle - |1\rangle$ . An extreme version of this situation is contained in the following proposition:

**Proposition 7** *Suppose that there are two states  $|0\rangle, |1\rangle$  whose corresponding density matrices  $\rho_{|0\rangle}, \rho_{|1\rangle}$  are perfectly distinguishable. Then the scheme is not an  $\epsilon$ -secure QAS for any  $\epsilon < 1$ .*

*Proof:* Since  $\rho_{|0\rangle}, \rho_{|1\rangle}$  can be distinguished, they must have orthogonal support, say on subspaces  $V_0, V_1$ . So consider an adversary who applies a phaseshift of  $-1$  conditioned on being in  $V_1$ . Then for all  $k$ ,  $\rho_{|0\rangle+|1\rangle}^{(k)}$  becomes  $\rho_{|0\rangle-|1\rangle}^{(k)}$ . Thus, Bob will decode the (orthogonal) state  $|0\rangle - |1\rangle$ .  $\square$

However, in general, the adversary cannot exactly distinguish two states, so we must allow some probability of failure. Note that it is sufficient in general to consider two encoded pure states, since any two mixed states can be written as ensembles of pure states, and the mixed states are distinguishable only if some pair of pure states are. Furthermore, we might as well let the two pure states be orthogonal, since if two nonorthogonal states  $|\psi_0\rangle$  and  $|\psi_1\rangle$  are distinguishable, two basis states  $|0\rangle$  and  $|1\rangle$  for the space spanned by  $|\psi_0\rangle$  and  $|\psi_1\rangle$  are at least as distinguishable.

Given the space limitations of this abstract, we outline the proof with a sequence of lemmas, whose proofs are contained in Appendix F.

We first consider the case when  $|0\rangle$  and  $|1\rangle$  can *almost* perfectly be distinguished. In that case, the adversary can change  $|0\rangle + |1\rangle$  to  $|0\rangle - |1\rangle$  with high (but not perfect) fidelity (stated formally in Lemma 16). When  $|0\rangle$  and  $|1\rangle$  are more similar, we first magnify the difference between them by repeatedly encoding the same state in multiple copies of the authentication scheme, then apply the above argument.

**Lemma 8** *Suppose that there are two states  $|0\rangle, |1\rangle$  such that  $D(\rho_{|0\rangle}, \rho_{|1\rangle}) \geq 1 - \eta$ . Then the scheme is not  $\epsilon$ -secure for  $|\psi\rangle = |0\rangle + |1\rangle$  for any  $\epsilon < 1 - 2\eta$ .*

When two states can be distinguished, but only just barely, the above lemma is not sufficient. Instead, we must magnify the distinguishability of the states  $|0\rangle$  and  $|1\rangle$  by repeating them by considering the tensor product of many copies of the same state. The probability of distinguishing then goes to 1 exponentially fast in the number of copies:

**Lemma 9** *Let  $\rho_0, \rho_1$  be density matrices with  $D(\rho_0, \rho_1) = \delta$ . Then  $D(\rho_0^{\otimes t}, \rho_1^{\otimes t}) \geq 1 - 2\exp(-t\delta^2/2)$ .*

We create these repeated states by encoding them in an iterated QAS consisting of  $t$  copies of the original QAS (with independent values of the key for each copy).

**Lemma 10** *Suppose we iterate the scheme  $t$  times. Let  $|\psi\rangle = \frac{1}{\sqrt{2}}(|000\dots 0\rangle + |111\dots 1\rangle)$ . If  $(A, B, \mathcal{K})$  is an  $\epsilon$ -secure QAS, then the iterated scheme is  $10t^3\epsilon$ -secure for the state  $|\psi\rangle$ .*

Note that the proof of this lemma goes through the following crucial claim, which follows from a simple hybrid argument.

**Claim 11 (Product states)** *The iterated scheme is  $t\epsilon$ -secure for any product state.*

Putting the various lemmas together, we find that, given two states  $|0\rangle$  and  $|1\rangle$  which are slightly distinguishable by the adversary, so  $D(\rho_0, \rho_1) \geq \delta$ , then in the iterated scheme,  $|000\dots 0\rangle$  and  $|111\dots 1\rangle$  are more distinguishable:  $D(\rho_{|000\dots 0\rangle}, \rho_{|111\dots 1\rangle}) \geq 1 - \eta$ , where  $\eta \leq 2 \exp(-t\delta^2/2)$ . Since the iterated scheme is  $10t^3\epsilon$ -secure for the state  $|\psi\rangle = \frac{1}{\sqrt{2}}(|000\dots 0\rangle + |111\dots 1\rangle)$ , then by the first lemma,

$$10t^3\epsilon > 1 - 2\eta \geq 1 - 4 \exp(-t\delta^2/2)$$

Choosing  $t = 1/\sqrt[3]{20\epsilon}$ , we get  $\delta \leq 4\epsilon^{1/6}$ .

## 7 Quantum Signatures

One consequence of the previous theorem is that digitally signing quantum messages is impossible. One can imagine more than one way of defining this task, but any reasonable definition must allow a recipient—who should not be able to alter signed messages—to learn something about the contents of the message. However, this is precisely what is forbidden by the previous theorem: in an information-theoretic setting, any adversary who can gain a non-trivial amount of information must be able to modify the authenticated state with non-negligible success.

If we consider computationally secure schemes, a somewhat narrower definition of digitally signing quantum states remains impossible to realize. If we assume a quantum digital signature protocol should allow any recipient to efficiently extract the original message, then a simple argument shows that he can also efficiently change it without being detected, contradicting the security of the scheme. Namely: Assume that there is transformation  $U$  with a small circuit which extracts the original message  $\rho$ , leaving auxiliary state  $|\varphi\rangle$  (which may not all be held by Bob). In order to preserve any entanglement between  $\rho$  and a reference system, the auxiliary state  $|\varphi\rangle$  must be independent of  $\rho$ . Therefore, Bob can replace  $\rho$  with any other state  $\rho'$  and then perform  $U^\dagger$  on  $\rho'$  and his portion of  $|\varphi\rangle$ , producing a valid signature for  $\rho'$ . This is an efficient procedure: the circuit for  $U^\dagger$  is just the circuit for  $U$  executed backwards.

Note that we have actually shown a somewhat stronger result: it is not possible, even when the sender is known to be honest, to authenticate a quantum message to a group of receivers (some of whom may be dishonest). This presentation also makes some limitations of our proof clear. For instance, the proof does not apply if the sender knows the identity of the quantum state he is signing, nor does it apply to signing classical messages.

## 8 Discussion and Conclusion

An interesting feature of our scheme: if the transmission quantum channel is not error free, we can modify our scheme to take advantage of the error-correction capability of the quantum code. More precisely, if  $\mathcal{B}$  rejects only when the number of observed errors is too large then error correction will fix natural noise or tampering of small amplitude.

We have examined various aspects of the problem of authenticating quantum messages. We have shown the security of a large class of private-key quantum authentication schemes, and presented a particular highly efficient scheme from that class. One feature of the scheme is that it completely encrypts the message, and we show that this is a necessary feature of any quantum message authentication code: if any observer can learn a substantial amount of information about the authenticated state, that observer also has a good chance of successfully changing the state without being detected. We have also studied authentication of quantum states in a public key context, and shown that while authentication is possible with public keys, digitally signing quantum states is never possible, even when only computational security is required.

The necessity for encryption is rather surprising, given that classical messages can be authenticated without encrypting them. The difference can be understood as a complementarity feature of quantum mechanics: authenticating a message in one basis requires encrypting it in the complementary Fourier-transformed basis. This is essentially another realization of the principle that measuring data in one basis disturbs it in any complementary basis. For classical messages, therefore, encryption is not required: only one basis is relevant. In contrast, for quantum messages, we require authentication in all bases and therefore we must also require encryption in all bases.

Note that purity-testing codes have many applications beyond QAS. For instance, the efficient purity-testing code of section 4 can be used to create a correspondingly efficient QKD protocol.

## Acknowledgments

We would like to thank Herbert Bernstein, Aart Blokhuis, Hoi Fung Chau, David DiVincenzo, Manny Knill, Debbie Leung, Michele Mosca, Eric Rains, and Ronald de Wolf for helpful discussions or comments.

## References

- [1] A. Ambainis, M. Mosca, A. Tapp and R. de Wolf “Private quantum channels”, *Proceedings of the 41st Annual Symposium on Foundations of Computer Science*, pp. 547 – 553, IEEE Computer Society Press, 2000.
- [2] A. Ashikhmin and E. Knill, “Nonbinary quantum stabilizer codes”, e-print quant-ph/0005008.
- [3] A. Ambainis, A. Smith and K. Yang, “Extracting Quantum Entanglement (General Entanglement Purification Protocols)”, *Proceedings of CCC 2002*, Montreal, May 2002.
- [4] C.H. Bennett, “Quantum Cryptography using any two nonorthogonal states”, *Phys. Rev. Lett.*, vol. 68, 1992, 3121 – 3124.
- [5] C. Bennett and G. Brassard. “Quantum cryptography: Public-key distribution and coin-tossing.” In *Proc. of IEEE Conference on Computers, Systems and Signal Processing*, 1984, pp. 175 – 179.
- [6] C. Bennett, G. Brassard, C. Crepeau, R. Jozsa, A. Peres and W. Wootters, “Teleporting an unknown quantum state via dual classical and EPR channels”, *Phys. Rev. Lett.*, pp. 1895 – 1899, March 1993.
- [7] C.H. Bennett, G. Brassard, S. Popescu, B. Schumacher, J. A. Smolin, W. K. Wootters, “Purification of Noisy Entanglement and Faithful Teleportation via Noisy Channels”, *Phys. Rev. Lett.*, vol. 76, 1996, pp. 722 – 725.
- [8] A. Beutelspacher and U. Rosenbaum, *Projective Geometry*, Cambridge University Press, Cambridge, 1998.
- [9] P. O. Boykin, V. Roychowdhury, “Optimal Encryption of Quantum Bits,” quant-ph/0003059.
- [10] A. R. Calderbank, E. M. Rains, P. W. Shor, and N. J. A. Sloane, “Quantum error correction via codes over  $GF(4)$ ”, *IEEE Trans. Inform. Theory*, vol. 44, 1998, pp. 1369 – 1387; quant-ph/9605005.
- [11] A. R. Calderbank and P. W. Shor, “Good Quantum Error-Correcting Codes Exist”, *Phys. Rev. A*, Vol. 54, No. 2, 1996, pp. 1098 – 1106.
- [12] J.L. Carter and M. N. Wegman, “New hash functions and their use in authentication and set equality”, *Journal of Computer and System Sciences*, Vol. 22, 1981, pp. 265 – 279.
- [13] A. K. Ekert. “Quantum cryptography based on Bell’s theorem”, *Physical Review Letters*, vol.67 pp.661 - 663, (1991).
- [14] D. Gottesman and I. Chuang, “Quantum digital signatures”, e-print quant-ph/0105032.

- [15] H.-K. Lo and H.F. Chau, “Unconditional Security Of Quantum Key Distribution Over Arbitrarily Long Distances”, *Science*, vol. 283, 1999, pp. 2050–2056.
- [16] M. A. Nielsen and I. L. Chuang, “Quantum Computation and Quantum Information”, *Cambridge University Press*, 676 pages, 2000.
- [17] T. Okamoto, K. Tanaka, S. Uchiyama. “Quantum public-key cryptosystems”, in *Proc. of CRYPTO 2000*, pp. 147-165.
- [18] A. Steane, “Multiple Particle Interference and Quantum Error Correction”, *Proc. Roy. Soc. Lond.*, A452, 1996, pp. 2551–2577
- [19] P.W. Shor and J. Preskill, “Simple Proof of Security of the BB84 Quantum Key Distribution Protocol”, *Phys. Rev. Lett.*, vol. 85, 2000, pp. 441–444.
- [20] S. Wiesner, “Conjugate coding”, *Sigact News*, Vol. 15, no. 1, 1983, pp. 78–88. Original manuscript written circa 1970.

## A Quantum Stabilizer Codes

A quantum error-correcting code (QECC) is a way of encoding quantum data (say  $m$  qubits) into  $n$  qubits ( $m < n$ ) such that the encoded data is protected from errors of small weight: the code is said to correct  $t$  errors if any operator which affects less than  $t$  qubits of the encoding can be corrected without disturbing the encoded state. Usually the goal in the construction of codes is to maximize this minimum distance for particular  $m, n$ . However, in this paper, we use the theory developed for those purposes to construct families of codes with a different type of property. For now, we review the necessary theory on a very general class of codes known as *stabilizer codes*.

Our construction is based on a class of QECCs for  $q$ -dimensional registers, with  $q = p^n$  a prime power (later we will specialize to the case where  $p = 2$ , so each register consists of  $n$  qubits). A basis for the set of all operators on the  $p$ -dimensional Hilbert space is the “shift/phase” error basis on  $p$ -dimensional Hilbert space, defined via  $E_{ab} = X^a Z^b$ , where  $\langle i|X|j \rangle = \delta_{i,j+1}$ ,  $\langle i|Z|j \rangle = \xi^i \delta_{i,j}$ , for  $\xi = \exp(2\pi i/p)$  a primitive  $p$ th root of unity, are the standard-basis matrix elements of the “shift by one” and “ramp the phase by one” operators. (Here, indices are in  $\mathbb{Z}_p$ .) This basis has a simple multiplication rule:  $E_{ab}E_{a'b'} = \xi^{a'b} E_{a+a', b+b'}$ . Thus,  $\{\xi^c E_{ab}\}$  is a group containing a basis for the whole operator space for one register. If we have  $n$  registers, we can simply use the tensor product  $E$  of  $n$  copies of this operator group; each element corresponds to a  $2n$ -dimensional vector, and the vectors  $x = (\mathbf{a}|\mathbf{b})$ ,  $y = (\mathbf{a}'|\mathbf{b}')$  come from commuting operators iff their symplectic inner product is 0 in  $\mathbb{Z}_p$ :

$$E_x E_y = E_y E_x \iff B(x, y) = \mathbf{a}' \cdot \mathbf{b} - \mathbf{a} \cdot \mathbf{b}' = 0. \quad (4)$$

A *stabilizer code* is a QECC given by an Abelian subgroup  $S$  of  $E$ , which does not contain any multiples of the identity other than  $I$  itself.  $S$  can be described by the set of  $2n$ -dimensional vectors  $x$  such that  $E_x \in S$ . This will be a subspace of  $\mathbb{Z}_p^{2n}$ . Moreover, it will be *totally isotropic*, i.e.  $B(x, y) = 0$  for all  $x, y$  in the subspace. If we take a set of generators for  $S$ , we can divide Hilbert space into a set of equidimensional orthogonal subspaces. Each such space  $T$  consists of common eigenvectors of all operators of  $S$  having a fixed pattern of eigenvalues, unique to  $T$ . The space with all eigenvalues  $+1$  is the “code space,” its elements are “codewords,” and the orthogonal spaces are labelled by “syndromes.”

Note that one can also view  $B(\cdot, \cdot)$  as a symplectic form over  $GF(p^{2n})$ , by choosing a set of generators for  $GF(p^{2n})$  as a vector space over  $\mathbb{Z}_p$ . By choosing different sets of generators for  $GF(p^{2n})$  as a vector space over  $\mathbb{Z}_p$ , we can get different symplectic forms  $B(\cdot, \cdot)$  over this finite vector space. By judicious choice of the generators, one can make  $B(\cdot, \cdot)$  correspond to *any* non-degenerate symplectic form over  $GF(p^{2n})$ .

**Undetectable errors** We can classify errors which lie in  $E$  into three categories: The errors corresponding to elements of  $Q$  are not truly errors—they leave the codewords unchanged. Errors which fail to commute with some element of  $Q$  move codewords into a subspace orthogonal to the code, so can be detected by the QECC. The remaining errors, those which commute with all elements in  $S$  but are not themselves in  $S$ , are the undetectable errors of the code. Thus, if  $Q^\perp$  is the space of vectors  $y$  for which  $B(x, y) = 0$  for all  $x \in Q$ , the set of undetectable errors is just  $Q^\perp - Q$ .

**Syndromes** Note that specifying the subgroup  $S$  by a set  $Q$  of elements of  $GF(p^{2n})$  isn't quite enough: operators differing by a phase  $\xi^c$  correspond to the same field element, but yield different QECC's in the Hilbert space. Given an  $s$ -dimensional totally isotropic subspace of  $\mathbb{Z}_p^{2n}$ , there are  $p^s$  possible choices of phases for the group  $S$ , which produce  $p^s$  different QECCs. However, all these codes have identical error correction properties. The corresponding code subspaces are all orthogonal and of the same dimension  $p^{n-s}$ . These codes are known as *cosets* of the code  $S$ , defined as the standard choice with all phases equal to 1.<sup>2</sup> The choice of phases is known as the *syndrome* (because errors outside  $S^\perp$  map the code into a different coset, and the syndrome therefore gives information about which error occurred). Measuring the syndrome projects a quantum state into one of these codes.

## B Alternative Security Definition

The definition of security of an authentication scheme given in Section 3 appears at first sight to have a major shortcoming: it does not tell what happens when  $\mathcal{A}$ 's input is a mixed state. Intuitively, this should not be a problem, since one expects security to extend from pure states to mixed states more or less by linearity. Indeed, this is the case, but it is not entirely clear what is *meant* by security when  $\mathcal{A}$ 's input is a mixed state  $\rho$ . One straightforward approach is to add a reference system  $R$ , and to assume the joint system of  $\mathcal{A}$  and  $R$  is always pure; then the requirement is that the final state of  $\mathcal{B}$  and  $R$  should high fidelity to the initial state. We could also use the following informal definition, which we will show is implied by Definition 2: as long as  $\mathcal{B}$ 's probability of acceptance is significant, then when he accepts, the fidelity of the message state he outputs to  $\mathcal{A}$ 's original state should be almost 1.

**Proposition 12** *Suppose that  $(A, B, \mathcal{K})$  is a  $\epsilon$ -secure QAS. Let  $\rho$  be the density matrix of  $\mathcal{A}$ 's input state and let  $\rho'$  be the density matrix output by  $\mathcal{B}$  conditioned on accepting the transmission as valid. Then if  $\mathcal{B}$ 's probability of accepting is  $p_{acc}$ , the fidelity of  $\rho$  to  $\rho'$  is bounded below. For any  $\rho$  and any adversary action  $\mathcal{O}$ , we have:  $F(\rho, \rho') \geq 1 - \frac{\epsilon}{p_{acc}}$ .*

In particular, if  $\epsilon$  is negligible and  $p_{acc}$  is non-negligible, then the fidelity of  $\mathcal{B}$ 's state to  $\mathcal{A}$ 's input state will be essentially 1.

To prove this, we first restate Proposition 12 more formally. Let  $\rho_{Bob}$  be the state of  $\mathcal{A}$ 's two output systems  $M, V$  when  $\mathcal{A}$ 's input is  $\rho$ . Denote the projector onto the space of accepting states by  $\Pi$ , that is  $\Pi = I_M \otimes |\text{ACC}\rangle\langle\text{ACC}|$ .

Using this notation,  $\mathcal{B}$ 's probability of accepting is  $p_{acc} = \text{Tr}(\Pi\rho_{Bob})$ , and the density matrix of the joint system  $M, V$  conditioned on acceptance is  $\rho_{acc} = \frac{\Pi\rho_{Bob}\Pi}{\text{Tr}(\Pi\rho_{Bob})} = \frac{\Pi\rho_{Bob}\Pi}{p_{acc}}$ .

Now since  $\rho_{acc}$  has been restricted to the cases where  $\mathcal{B}$  accepts, we can write  $\rho_{acc} = \rho' \otimes |\text{ACC}\rangle\langle\text{ACC}|$ , where  $\rho'$  is the density matrix of  $\mathcal{B}$ 's message system conditioned on his acceptance of the transmission as valid. From the definition of fidelity, we can see that

$$F(\rho, \rho') = F(\rho \otimes |\text{ACC}\rangle\langle\text{ACC}|, \rho_{acc})$$

<sup>2</sup>Actually, the "standard" coset also depends on the selection of a basis of generators for  $S$ .

We can now restate the theorem:

**Claim (Proposition 12):**  $F(\rho \otimes |\text{ACC}\rangle\langle\text{ACC}|, \rho_{acc}) \geq 1 - \frac{\epsilon}{p_{acc}}$

*Proof (of Theorem 12):* Write  $\rho = \sum_i p_i |\psi_i\rangle\langle\psi_i|$  for some orthonormal basis  $\{|\psi_i\rangle\}$ . For each  $i$ , let  $\rho_i$  be  $\mathcal{B}$ 's output when  $\mathcal{A}$  uses input  $|\psi_i\rangle$ . We have  $\rho_{Bob} = \sum_i p_i \rho_i$ .

For each  $i$ , let  $P_i = |\psi_i\rangle\langle\psi_i| \otimes |\text{ACC}\rangle\langle\text{ACC}|$  and let  $Q_i = (I_M - |\psi_i\rangle\langle\psi_i|) \otimes |\text{ACC}\rangle\langle\text{ACC}|$  so that  $P_i + Q_i = \Pi$ .

Now we can write  $\rho \otimes |\text{ACC}\rangle\langle\text{ACC}| = \sum_i p_i P_i$ , and  $\rho_{acc} = \sum_i p_i \frac{\Pi \rho_i \Pi}{p_{acc}}$ . By the concavity of fidelity (Theorem 9.7 of [16]), we get

$$F(\rho \otimes |\text{ACC}\rangle\langle\text{ACC}|, \rho_{acc}) = F\left(\sum_i p_i P_i, \sum_i \frac{\Pi \rho_i \Pi}{p_{acc}}\right) \geq \sum_i p_i F\left(P_i, \frac{\Pi \rho_i \Pi}{p_{acc}}\right) \quad (5)$$

The formula for fidelity for one-dimensional projectors is simple: for a projector  $P$  and any density matrix  $\sigma$ , we have  $F(P, \sigma) = \sqrt{\text{Tr}(P\sigma)}$ . Thus expression (5) simplifies to

$$\sum_i p_i \sqrt{\text{Tr}\left(P_i \frac{\Pi \rho_i \Pi}{p_{acc}}\right)}$$

Using the fact that  $\Pi P_i \Pi = P_i$ , we can further simplify this:

$$\sum_i p_i \sqrt{\frac{\text{Tr}(P_i \rho_i)}{p_{acc}}}$$

Since  $\frac{\text{Tr}(P_i \rho_i)}{p_{acc}}$  is always less than 1, we can obtain a lower bound by removing the square root sign:

$$F(\rho \otimes |\text{ACC}\rangle\langle\text{ACC}|, \rho_{acc}) \geq \frac{\sum_i p_i \text{Tr}(P_i \rho_i)}{p_{acc}} \quad (6)$$

Now the acceptance probability  $p_{acc} = \text{Tr}(\Pi \rho_{Bob})$  can be written as  $\sum_i p_i \text{Tr}(\Pi \rho_i)$ . Using the fact that  $\Pi = P_i + Q_i$  we get that  $p_{acc} = (\sum_i p_i \text{Tr}(P_i \rho_i)) + (\sum_i p_i \text{Tr}(Q_i \rho_i))$ .

But by the definition of  $\epsilon$ -security, we know that for each  $i$ , we have  $\text{Tr}(Q_i \rho_i) \leq \epsilon$ , and so  $p_{acc} \leq (\sum_i p_i \text{Tr}(P_i \rho_i)) + \epsilon$ , and so we get  $(\sum_i p_i \text{Tr}(P_i \rho_i)) \geq p_{acc} - \epsilon$ . Applying this observation to expression (6), we get :

$$F(\rho \otimes |\text{ACC}\rangle\langle\text{ACC}|, \rho_{acc}) \geq \frac{p_{acc} - \epsilon}{p_{acc}} = 1 - \frac{\epsilon}{p_{acc}}$$

□

## C Proof of Proposition 1

Proposition 1 states that a stabilizer purity testing code can always be used to produce a purity testing protocol with the same error  $\epsilon$ .

*Proof:* If  $\mathcal{A}$  and  $\mathcal{B}$  are given  $n$  EPR pairs, this procedure will always accept, and the output will always be  $|\Phi^+\rangle^{\otimes n}$ . Thus,  $\mathcal{T}$  satisfies the completeness condition.

Suppose for the moment that the input state is  $(E_x \otimes I)|\Phi^+\rangle^{\otimes n}$ , for  $E_x \in E$ ,  $x \neq 0$ . Then when  $k$  is chosen at random, there is only probability  $\epsilon$  that  $x \in Q_k^\perp - Q_k$ . If  $x \notin Q_k^\perp$ , then  $\mathcal{A}$  and  $\mathcal{B}$  will find different error syndromes, and therefore reject the state. If  $x \in Q_k^\perp$ , then  $\mathcal{A}$  and  $\mathcal{B}$  will accept the state, but if  $x \in Q_k$ ,

then the output state will be  $|\Phi^+\rangle^{\otimes m}$  anyway. Thus, the probability that  $\mathcal{A}$  and  $\mathcal{B}$  will accept an incorrect state is at most  $\epsilon$ .

To prove the soundness condition, we can use this fact and a technique of Lo and Chau [15]. The states  $(E_x \otimes I)|\Phi^+\rangle^{\otimes n}$  form the Bell basis for the Hilbert space of  $\mathcal{A}$  and  $\mathcal{B}$ . Suppose a nonlocal third party first measured the input state  $\rho$  in the Bell basis; call this measurement  $B$ . Then the argument of the previous paragraph would apply to show the soundness condition. In fact, it would be sufficient if Alice and Bob used the nonlocal measurement  $Q_k \otimes Q_k$  which compares the  $Q_k$ -syndromes for  $\mathcal{A}$  and  $\mathcal{B}$  without measuring them precisely. This is a submeasurement of the Bell measurement  $B$  — that is, it gives no additional information about the state. Therefore it commutes with  $B$ , so the sequence  $B$  followed by  $Q_k \otimes Q_k$  is the same as  $Q_k \otimes Q_k$  followed by  $B$ , which therefore gives probability at least  $1 - \epsilon$  of success for general input states  $\rho$ . But if the state after  $Q_k \otimes Q_k$  gives, from a Bell measurement,  $|\Phi^+\rangle^{\otimes m}$  or  $|\text{REJ}\rangle$  with probability  $1 - \epsilon$ , then the state itself must have fidelity  $1 - \epsilon$  to the projection  $P$ . Therefore, the measurement  $Q_k \otimes Q_k$  without  $B$  satisfies the soundness condition. Moreover,  $\mathcal{A}$  and  $\mathcal{B}$ 's actual procedure  $\mathcal{T}$  is a refinement of  $Q_k \otimes Q_k$ —that is, it gathers strictly more information. Therefore, it also satisfies the soundness condition, and  $\mathcal{T}$  is a purity testing protocol with error  $\epsilon$ . □

## D Analysis of Purity-Testing Code Construction

It is straightforward to extend the purity testing code defined in Section 4.1 to arbitrary finite fields  $GF(q)$ . To do so, we work over a global field  $GF(q^{2rs})$  and break it down into both a  $2r$ -dimensional vector space over  $GF(q^s)$  and a  $2rs$ -dimensional vector space over  $GF(q)$ . We exploit this by defining our  $GF(2)$ -valued symplectic form  $B$  via a choice of a  $GF(q^s)$ -valued symplectic form  $C$  on  $GF(q^{2rs})$  and a non-null linear map  $L : GF(q^s) \rightarrow GF(q)$ , where linearity is defined by viewing  $GF(q^s)$  as an  $s$ -dimensional vector space over  $GF(q)$ . Then

$$B(x, y) := L(C(x, y)) . \quad (7)$$

Bilinearity and alternation of  $B$  are obvious. For fixed  $y(x)$ , by  $C$ 's nondegeneracy there is a  $z$  such that  $C(z, y)C(x, z) \neq 0$ . Considering  $\alpha z$  in place of  $z$ , for all scalars  $\alpha \in GF(2^s)$ , and still holding  $y(x)$  fixed, shows (by bilinearity of  $C$ ) that  $C(x, y)$  takes all values in  $GF(q^s)$  as  $x(y)$  is varied; by non-nullity of  $L$ , not all of these can map to zero, i.e.  $B$  is nondegenerate.

The definition of the purity testing code  $\{Q_k\}$  is then the same as in the binary case.

**Theorem 13** *The set of codes  $Q_k$  form a stabilizer purity testing code with error*

$$\epsilon = \frac{2r}{q^s + 1} . \quad (8)$$

*Each code  $Q_k$  encodes  $m = (r - 1)s$  dimension  $q$  registers in  $n = rs$  registers.*

We must show (a) that  $Q_k$  is totally isotropic, and (b) that the error probability is at most  $\epsilon$ .

(a) For  $\alpha, \beta \in GF(2^s)$ , we have

$$B(\alpha x, \beta y) = L(C(\alpha x, \beta y)) = L(\alpha\beta C(x, y)) . \quad (9)$$

Fix an  $x \in Q_k - \{0\}$ . Every  $y \in Q_k$  may be written as  $\alpha x$ , for some  $\alpha \in GF(q^s)$ . Now  $C(x, x) = 0$ . So,  $B(\alpha x, \beta x) = 0$ , i.e.,  $Q_k$  is totally isotropic under the symplectic form  $B$ .



(b) We must find, for an arbitrary error  $E_x$  (which can be described via a  $2n$ -dimensional  $GF(q)$  vector  $x$ ), an upper bound on the number of  $Q_k^\perp - Q_k$  it can belong to. It will be sufficient to bound the number of  $Q_k^\perp$  the error can belong to, since  $|Q_k|$  is small compared to  $|Q_k^\perp|$  in our context.  $x \in Q_k^\perp$  means  $B(x, y) = 0$  for all  $y \in Q_k$ . By choice of  $s$  linearly independent  $y \in Q_k$  this imposes  $s$  linearly independent linear equations on  $x$ . We will show below that if we take any  $2r$  codes  $Q_k$  defined by points on  $\Upsilon$ , and take  $s$  independent vectors from each, the resulting set of  $2rs$  vectors is linearly independent. Thus if  $E_x$  is undetectable in  $2r$  such codes, this imposes the dimension's worth ( $2rs$ ) of linearly independent equations on  $x$ . Consequently,  $E_x$  must be detectable in all the remaining codes, i.e.,  $E_x$  can satisfy  $x \in Q_k^\perp$  for *at most*  $2r$  values of  $k$ , when  $Q_k$  are chosen among the  $q^s + 1$  available  $s$ -dimensional spaces corresponding to points on  $\Upsilon$ . Thus, the  $\{Q_k\}$  form a purity testing code with error

$$\epsilon \leq \frac{2r}{q^s + 1}. \quad (10)$$

We now show the claimed property of codes defined by  $\Upsilon$ . A set of points in a projective geometry of dimension  $d-1$  are said to be in general position if any  $d$  (= dimension of the underlying vector space, when, as in our case, such exists) of them are linearly independent. The points on the normal rational curve  $\Upsilon$  are in general position, and indeed a maximal set of such points. (To verify that they are in general position one shows that for any  $2r$  points on the curve, the determinant of the matrix of their coordinates is nonzero; these are Vandermonde determinants.) That is, any  $2r$  points on  $\Upsilon$  are linearly independent. Each point  $k$  on  $\Upsilon$  corresponds to an  $s$ -dimensional code  $Q_k$ , consisting of  $2rs$ -dimensional vectors. Let  $z$  be any nonzero element of  $Q_k$ . As  $\alpha$  ranges over  $GF(q^s)$ ,  $\alpha z$  ranges over all vectors in  $Q_k$ . Thus, if any vector from  $Q_k$  is a linear combination of vectors from other codes  $\{Q_j\}$ , then all of  $Q_k$  is also a linear combination of vectors from  $\{Q_j\}$ , and  $k$  is linearly dependent on the points  $\{j\}$  of  $\Upsilon$ . So if we take any  $2r$  codes  $Q_k$ , and take  $s$  independent vectors from each, the resulting set of  $2rs$  vectors is linearly independent.

## E Proof of secure authentication

Corollary 3 states that the interactive authentication protocol 5.1 is secure.

*Proof (of Corollary 3):*

The completeness of the protocol can be seen by inspection: in the absence of intervention,  $\mathcal{A}$  and  $\mathcal{B}$  will share the Bell states  $|\Phi^+\rangle^{\otimes m}$  at the end of step 6 and so after the teleportation in step 7  $\mathcal{B}$ 's output will be exactly the input of  $\mathcal{A}$ .

To prove soundness, suppose that  $\mathcal{A}$ 's input is a pure state  $|\psi\rangle$ . Intuitively, at the end of step 6,  $\mathcal{A}$  and  $\mathcal{B}$  share something very close to  $|\Phi^+\rangle^{\otimes m}$ , and so after the teleportation in step 7 either  $\mathcal{B}$ 's output will be very close to  $\mathcal{A}$ 's input, or he will reject because of interference from the adversary.

More formally, after step 6, the joint state  $\rho_{AB}$  satisfies  $\text{Tr}(P\rho_{AB}) \geq 1 - \epsilon$ . At this point, by assumption the only thing that the adversary can do is attempt to jam the communication between  $\mathcal{A}$  and  $\mathcal{B}$ . Thus the effect of step 7 will be to map the subspace given by  $P$  into the subspace given by  $P_1^{|\psi\rangle}$ . Consequently, at the end of the protocol,  $\mathcal{B}$ 's output density matrix will indeed lie almost completely in the subspace defined by  $P_1^{|\psi\rangle}$ .

□

Theorem 4 states that the non-interactive Protocol 5.2 is secure. To prove this, we show that Protocol 5.1 is equivalent to 5.2, by moving through two intermediate protocols E.1 and E.2. We reduce the security of each protocol to the previous one; since Protocol 5.1 is secure by Corollary 3, the theorem follows.

PROTOCOL 5.1  $\rightarrow$  PROTOCOL E.1: We obtain protocol E.1 by observing that in protocol 5.1,  $\mathcal{A}$  can perform all of her operations (except for the transmissions) *before* she actually sends anything to  $\mathcal{B}$ , since

**Protocol E.1 ( Intermediate Protocol I )**

- 1:**  $\mathcal{A}$  and  $\mathcal{B}$  agree on some stabilizer purity testing code  $\{Q_k\}$
- 2:**  $\mathcal{A}$  generates  $2n$  qubits in state  $|\Phi^+\rangle^{\otimes n}$ .  $\mathcal{A}$  picks at random  $k \in \mathcal{K}$ , and measures the syndrome  $y$  of the stabilizer code  $Q_k$  on the first half of the EPR pairs.  $\mathcal{A}$  decodes her  $n$ -qubit word according to  $Q_k$ .  $\mathcal{A}$  performs the Bell measurement to start teleportation with her state  $\rho$ , using the decoded state as if it were half of  $|\Phi^+\rangle$  pairs, but does not yet reveal the measurement results  $x$  of the teleportation.  $\mathcal{A}$  sends the second half of each EPR pair to  $\mathcal{B}$ .
- 3:**  $\mathcal{B}$  announces that he has received the  $n$  qubits. Denote the received state by  $\sigma'$ .
- 4:**  $\mathcal{A}$  announces  $k$  and the syndrome  $y$  of  $Q_k$  to  $\mathcal{B}$ .
- 5:**  $\mathcal{B}$  measures the syndrome  $y'$  of  $Q_k$  on his  $n$  qubits.  $\mathcal{B}$  compares the syndrome  $y'$  to  $y$ . If they are different,  $\mathcal{B}$  aborts.  $\mathcal{B}$  decodes his  $n$ -qubit word according to  $Q_k$ .
- 6:**  $\mathcal{A}$  concludes the teleportation by sending the teleportation measurement results  $x$  from step 2.  $\mathcal{B}$  does his part of the teleportation and obtains  $\rho'$ .

**Protocol E.2 ( Intermediate Protocol II )**

- 1:**  $\mathcal{A}$  and  $\mathcal{B}$  agree on some stabilizer purity testing code  $\{Q_k\}$
- 2:**  $\mathcal{A}$  chooses a random  $2n$  bit key  $x$  and  $q$ -encrypts  $\rho$  as  $\tau$  using  $x$ .  $\mathcal{A}$  picks a random  $k \in \mathcal{K}$  and syndrome  $s$  for the code  $Q_k$  and encodes  $\tau$  according to  $Q_k$ .  $\mathcal{A}$  sends the result to  $\mathcal{B}$ .
- 3:**  $\mathcal{B}$  announces that he has received the  $n$  qubits. Denote the received state by  $\sigma'$ .
- 4:**  $\mathcal{A}$  announces  $k$ ,  $x$ , and  $y$  to  $\mathcal{B}$ .
- 5:**  $\mathcal{B}$  measures the syndrome  $y'$  of the code  $Q_k$ .  $\mathcal{B}$  compares  $y$  to  $y'$ , and aborts if they are different.  $\mathcal{B}$  decodes his  $n$ -qubit word according to  $Q_k$ , obtaining  $\tau'$ .  $\mathcal{B}$   $q$ -decrypts  $\tau'$  using  $x$  and obtains  $\rho'$ .

these actions do not depend on  $\mathcal{B}$ 's feedback. This will not change any of the states transmitted in the protocol or computed by Bob, and so both completeness and soundness will remain the same.

PROTOCOL E.1  $\rightarrow$  PROTOCOL E.2: There are two changes between Protocols E.1 and E.2. First, note that measuring the first qubit of a state  $|\Phi^+\rangle$  and obtaining a random bit  $c_i$  is equivalent to choosing  $c_i$  at random and preparing the pure state  $|c_i\rangle \otimes |c_i\rangle$ . Therefore, instead of preparing the state  $|\Phi^+\rangle^{\otimes n}$  and measuring the syndrome of half of it,  $\mathcal{A}$  may as well choose the syndromes  $s$  at random and encode both halves of the state  $|\Phi^+\rangle^{\otimes m}$  using the code  $Q_k$  and the syndrome  $s$ .

Second, rather than teleporting her state  $\rho$  to  $\mathcal{B}$  using the EPR halves which were encoded in  $Q_{s_1, s_2}$ ,  $\mathcal{A}$  can encrypt  $\rho$  using a quantum one-time pad (QOTP) and send it to  $\mathcal{B}$  directly, further encoded in  $Q_k$ . These behaviours are equivalent since either way, the encoded state is  $\sigma_x^{\vec{t}_1} \sigma_z^{\vec{t}_2} \rho \sigma_z^{\vec{t}_2} \sigma_x^{\vec{t}_1}$ , where  $\vec{t}_1$  and  $\vec{t}_2$  are random  $n$ -bit vectors.

PROTOCOL E.2  $\rightarrow$  PROTOCOL 5.2: In Protocol 5.2, all the random choices of  $\mathcal{A}$  are replaced with the bits taken from a secret random key shared only by her and  $\mathcal{B}$ . This eliminates the need for an authenticated classical channel, and for any interaction in the protocol. This transformation can only increase the security of the protocol as it simply removes the adversary's ability to jam the classical communication.  $\square$

## F Proofs from Section 6

**Theorem 14 (Main Lower Bound)** A QAS with error  $\epsilon$  is an encryption scheme with error at most  $4\epsilon^{1/6}$ .

To get a sense of the proof, consider the following proposition:

**Proposition 15** Suppose that there are two states  $|0\rangle, |1\rangle$  whose corresponding density matrices  $\rho_{|0\rangle}, \rho_{|1\rangle}$  are perfectly distinguishable. Then the scheme is not an  $\epsilon$ -secure QAS for any  $\epsilon < 1$ .

*Proof:* Since  $\rho_{|0\rangle}, \rho_{|1\rangle}$  can be distinguished, they must have orthogonal support, say on subspaces  $V_0, V_1$ . So consider an adversary who applies a phaseshift of  $-1$  conditioned on being in  $V_1$ . Then for all  $k$ ,  $\rho_{|0\rangle+|1\rangle}^{(k)}$  becomes  $\rho_{|0\rangle-|1\rangle}^{(k)}$ . Thus, Bob will decode the (orthogonal) state  $|0\rangle - |1\rangle$ .  $\square$

However, in general, the adversary cannot exactly distinguish two states, so we must allow some probability of failure. Note that it is sufficient in general to consider two encoded pure states, since any two mixed states can be written as ensembles of pure states, and the mixed states are distinguishable only if some pair of pure states are. Furthermore, we might as well let the two pure states be orthogonal, since if two nonorthogonal states  $|\psi_0\rangle$  and  $|\psi_1\rangle$  are distinguishable, two basis states  $|0\rangle$  and  $|1\rangle$  for the space spanned by  $|\psi_0\rangle$  and  $|\psi_1\rangle$  are at least as distinguishable.

We first consider the case when  $|0\rangle$  and  $|1\rangle$  can *almost* perfectly be distinguished. In that case, the adversary can change  $|0\rangle + |1\rangle$  to  $|0\rangle - |1\rangle$  with high (but not perfect) fidelity (stated formally in Lemma 16). When  $|0\rangle$  and  $|1\rangle$  are more similar, we first magnify the difference between them by repeatedly encoding the same state in multiple copies of the authentication scheme, then apply the above argument.

**Lemma 16** Suppose that there are two states  $|0\rangle, |1\rangle$  such that  $D(\rho_{|0\rangle}, \rho_{|1\rangle}) \geq 1 - \eta$ . Then the scheme is not  $\epsilon$ -secure for  $|\psi\rangle = |0\rangle + |1\rangle$  for any  $\epsilon < 1 - 2\eta$ .

*Proof (of Lemma 16):* Let  $\rho_0 = \rho_{|0\rangle}$  and  $\rho_1 = \rho_{|1\rangle}$ . Consider the Hermitian matrix  $\sigma = \rho_0 - \rho_1$ . We can diagonalize  $\sigma$ . Let  $V_0$  be the space spanned by eigenvectors with non-negative eigenvalues and let  $V_1$  be the orthogonal complement.

Since  $1/2 \text{Tr}|\sigma| \geq 1 - \eta$ , but  $\text{Tr} \sigma = 0$ , we know that  $\text{Tr}(V_0\sigma) = -\text{Tr}(V_1\sigma) \geq 1 - \eta$ . Thus,  $\text{Tr}(V_0\rho_0) \geq \text{Tr}(V_0\sigma) \geq 1 - \eta$ . Similarly,  $\text{Tr}(V_1\rho_1) \geq -\text{Tr}(V_1\sigma) \geq 1 - \eta$ .

Consider an adversary who applies a phaseshift of  $-1$  conditioned on being in  $V_1$ . Fix a particular key  $k$ . Let  $p_0 = \text{Tr}(V_0\rho_0^{(k)})$  and  $p_1 = \text{Tr}(V_1\rho_1^{(k)})$ . We know that the expected values of  $p_0$  and  $p_1$  are both at least  $1 - \eta$ .

**Claim 17** When the input state is  $\frac{1}{\sqrt{2}}(|0\rangle + |1\rangle)$ , the fidelity of Bob's output to the state  $\frac{1}{\sqrt{2}}(|0\rangle - |1\rangle)|\text{ACC}\rangle$  is at least  $p_0 + p_1 - 1$ .

*Proof:* Consider some reference system  $R$  which allows us to purify the states  $\rho_0^{(k)}, \rho_1^{(k)}$  to the states  $|\tilde{0}\rangle, |\tilde{1}\rangle$ . Let  $|\tilde{\psi}\rangle$  be the image of  $\frac{1}{\sqrt{2}}(|\tilde{0}\rangle + |\tilde{1}\rangle)$  under the adversary's conditional phaseshift.

We want to show that  $|\tilde{\psi}\rangle$  is close to a correct encoding of  $\frac{1}{\sqrt{2}}(|0\rangle - |1\rangle)$ , i.e. close to

$$\frac{1}{\sqrt{2}}(|\tilde{0}\rangle - |\tilde{1}\rangle) = \frac{1}{\sqrt{2}}(V_0|\tilde{0}\rangle + V_1|\tilde{0}\rangle - V_0|\tilde{1}\rangle - V_1|\tilde{1}\rangle).$$

After the transformation, we obtain

$$|\tilde{\psi}\rangle = \frac{1}{\sqrt{2}}(V_0|\tilde{0}\rangle - V_1|\tilde{0}\rangle + V_0|\tilde{1}\rangle - V_1|\tilde{1}\rangle).$$

Thus,

$$\begin{aligned}
\langle \tilde{\psi} | \frac{1}{\sqrt{2}}(|\tilde{0}\rangle - |\tilde{1}\rangle) &= \frac{1}{2} (\langle \tilde{0} | V_0 | \tilde{0} \rangle - \langle \tilde{0} | V_1 | \tilde{0} \rangle - \langle \tilde{1} | V_0 | \tilde{1} \rangle + \langle \tilde{1} | V_1 | \tilde{1} \rangle \\
&\quad - \langle \tilde{0} | V_0 | \tilde{1} \rangle + \langle \tilde{1} | V_0 | \tilde{0} \rangle + \langle \tilde{0} | V_1 | \tilde{1} \rangle - \langle \tilde{1} | V_1 | \tilde{0} \rangle) \\
&= \frac{1}{2} \left( \text{Tr}(V_0 \rho_0^{(k)}) - \text{Tr}(V_1 \rho_0^{(k)}) - \text{Tr}(V_0 \rho_1^{(k)}) + \text{Tr}(V_1 \rho_1^{(k)}) \right. \\
&\quad \left. - [\langle \tilde{0} | V_0 | \tilde{1} \rangle - \langle \tilde{1} | V_0 | \tilde{0} \rangle] + [\langle \tilde{0} | V_1 | \tilde{1} \rangle - \langle \tilde{1} | V_1 | \tilde{0} \rangle] \right).
\end{aligned}$$

We can substitute for the first line in terms of  $p_0$  and  $p_1$ , which are real. The second line is purely imaginary. Thus,

$$\left| \langle \tilde{\psi} | \frac{1}{\sqrt{2}}(|\tilde{0}\rangle - |\tilde{1}\rangle) \right| \geq \frac{1}{2} [p_0 - (1 - p_0) - (1 - p_1) + p_1] = p_0 + p_1 - 1.$$

Bob's decoding can only increase the fidelity of the two states, as can discarding the reference system, proving the claim.  $\square$

Thus, for a specific value  $k$  of the key,  $F(\rho^{(k)}, \frac{1}{\sqrt{2}}(|0\rangle - |1\rangle)|_{\text{ACC}}) \geq p_0 + p_1 - 1$ , where  $\rho^{(k)}$  is the output after the adversary's transformation when the input is  $\frac{1}{\sqrt{2}}(|0\rangle + |1\rangle)$ . Fidelity is concave, so by Jensen's inequality the fidelity of the average density matrix  $\rho = \frac{1}{|\mathcal{K}|} \sum_k \rho^{(k)}$  is at least the average of the fidelities for each  $k$ . That is,

$$F(\rho, \frac{1}{\sqrt{2}}(|0\rangle - |1\rangle)|_{\text{ACC}}) \geq \frac{1}{|\mathcal{K}|} \sum_k (p_0 + p_1 - 1) \geq 1 - 2\eta.$$

In other words, the adversary can change the state  $\frac{1}{\sqrt{2}}(|0\rangle + |1\rangle)$  with probability at least  $1 - 2\eta$ .  $\square$

When two states can be distinguished, but only just barely, the above lemma is not sufficient. Instead, we must magnify the distinguishability of the states  $|0\rangle$  and  $|1\rangle$  by repeating them by considering the tensor product of many copies of the same state. The probability of distinguishing then goes to 1 exponentially fast in the number of copies:

**Lemma 18** *Let  $\rho_0, \rho_1$  be density matrices with  $D(\rho_0, \rho_1) = \delta$ . Then  $D(\rho_0^{\otimes t}, \rho_1^{\otimes t}) \geq 1 - 2 \exp(-t\delta^2/2)$ .*

*Proof (of Lemma 18):* We can bound  $D(\rho_0^{\otimes t}, \rho_1^{\otimes t})$  by giving a test which distinguishes them very well. We know there exists a measurement given by spaces  $V_0, V_1$  such that  $\text{Tr}(V_0 \rho_0) - \text{Tr}(V_0 \rho_1) = \delta$ . Consider the test which performs this measurement on  $t$  independent copies of  $\rho_0$  or  $\rho_1$ . The test outputs 0 if more than  $(\text{Tr}(V_0 \rho_0) + \text{Tr}(V_0 \rho_1))/2$  of the measurements produce 0.

By the Chernoff bound, the probability that this test will make the wrong guess is at most  $\exp(-t\delta^2/2)$ . Thus,  $D(\rho_0^{\otimes t}, \rho_1^{\otimes t}) \geq 1 - 2 \exp(-t\delta^2/2)$ .  $\square$

We create these repeated states by encoding them in an iterated QAS consisting of  $t$  copies of the original QAS (with independent values of the key for each copy).

**Lemma 19** *Suppose we iterate the scheme  $t$  times. Let  $|\psi\rangle = \frac{1}{\sqrt{2}}(|000\dots 0\rangle + |111\dots 1\rangle)$ . If  $(A, B, \mathcal{K})$  is an  $\epsilon$ -secure QAS, then the iterated scheme is  $10t^3\epsilon$ -secure for the state  $|\psi\rangle$ .*

Note that the proof of this lemma goes through the following crucial claim, which follows from a simple hybrid argument.

**Claim 20 (Product states)** *The iterated scheme is  $t\epsilon$ -secure for any product state.*

*Proof (of Claim 20):* For simplicity we prove the claim for the state  $|000\dots 0\rangle$ . The same proof works for any product pure state (and in fact for separable states in general).

Intuitively, an adversary who modifies the state  $|000\dots 0\rangle$  must change some component of the state. We can formalize this by rewriting the projector  $P_0^{000\dots 0}$  in terms of the individual projectors  $P_0^{0^i}$ .

For the case  $t = 2$ , Bob accepts only if he finds the verification qubits for both schemes in the accept state.

$$\begin{aligned} P_0^{00} &= (I_{m_1 m_2} - |00\rangle\langle 00|) \otimes |\text{ACC}_1\rangle\langle \text{ACC}_1| \otimes |\text{ACC}_2\rangle\langle \text{ACC}_2| \\ &= \left( (I_{m_1} - |0\rangle\langle 0|) \otimes I_{m_2} + I_{m_1} \otimes (I_{m_2} - |0\rangle\langle 0|) - (I_{m_1} - |0\rangle\langle 0|) \otimes (I_{m_2} - |0\rangle\langle 0|) \right) \\ &\quad \otimes |\text{ACC}_1\rangle\langle \text{ACC}_1| \otimes |\text{ACC}_2\rangle\langle \text{ACC}_2| \\ &= P_0^{0^1} \otimes |\text{ACC}_2\rangle\langle \text{ACC}_2| + P_0^{0^2} \otimes |\text{ACC}_1\rangle\langle \text{ACC}_1| - P_0^{0^1} \otimes P_0^{0^2} \end{aligned}$$

Since  $P_0^{0^1} \otimes P_0^{0^2}$  is positive, for all  $\rho$ , we have

$$\text{Tr}(P_0^{00} \rho) \leq \text{Tr}(P_0^{0^1} \rho) + \text{Tr}(P_0^{0^2} \rho) \leq 2\epsilon$$

Similarly, for larger values of  $t$  we have

$$\text{Tr}(P_0^{000\dots 0} \rho) \leq \sum_{i=1}^t \text{Tr}(P_0^{0^i} \rho) \leq t\epsilon$$

Thus the iterated scheme is  $t\epsilon$ -secure for  $|000\dots 0\rangle$  (and in fact for all product states).  $\square$

*Proof (of Lemma 19):* Consider the net superoperator due to encoding, decoding, and the adversary's intervention, i.e.  $\mathcal{O}_{net} = \frac{1}{|\mathcal{K}|} \sum_k B_k \mathcal{O}_{adv} A_k$ . By introducing an ancilla system  $R$ , we can extend this superoperator to a linear transformation on the joint system  $M \otimes R \otimes V$  (where  $M$  is the message system and  $V$  is Bob's verification qubit). For a pure state  $|\psi\rangle$ , write its image as

$$|\psi\rangle|\gamma_{|\psi}\rangle|\text{ACC}\rangle + |\beta_{|\psi}\rangle|\text{REJ}\rangle + |\delta_{|\psi}\rangle|\text{ACC}\rangle$$

where  $|\delta_{|\psi}\rangle$  is a joint state of  $MR$  which is orthogonal to the subspace  $|\psi\rangle \otimes R$ .

Now consider the family of states  $|\psi_i\rangle = |\underbrace{000\dots 0}_i \underbrace{111\dots 1}_{t-i}\rangle$ , and let  $|\gamma_i\rangle = |\gamma_{|\psi_i}\rangle$  and  $|\delta_i\rangle = |\delta_{|\psi_i}\rangle$ .

**Claim 21** *For all  $i = 0, \dots, t-1$ , we have  $\|\frac{1}{2}(|\gamma_{i+1}\rangle - |\gamma_i\rangle)\| \leq (1 + \sqrt{2})\sqrt{t\epsilon}$*

*Proof:* Fix  $i$ . Note that  $|\psi_+\rangle = \frac{1}{\sqrt{2}}(|\psi_{i+1}\rangle + |\psi_i\rangle)$  is a product state (with  $H|0\rangle$  in one position), as is  $|\psi_-\rangle = \frac{1}{\sqrt{2}}(|\psi_{i+1}\rangle - |\psi_i\rangle)$ . The image of  $|\psi_+\rangle$  can be written

$$\begin{aligned} &\frac{1}{\sqrt{2}} \left( (|\psi_{i+1}\rangle|\gamma_{i+1}\rangle + |\psi_i\rangle|\gamma_i\rangle)|\text{ACC}\rangle + (|\delta_{i+1}\rangle + |\delta_i\rangle)|\text{ACC}\rangle + (|\beta_{i+1}\rangle + |\beta_i\rangle)|\text{REJ}\rangle \right) \\ &= \left( |\psi_+\rangle \frac{1}{2}(|\gamma_{i+1}\rangle + |\gamma_i\rangle) + |\psi_-\rangle \frac{1}{2}(|\gamma_{i+1}\rangle - |\gamma_i\rangle) + \frac{1}{\sqrt{2}}(|\delta_{i+1}\rangle + |\delta_i\rangle) \right) |\text{ACC}\rangle \\ &\quad + \frac{1}{\sqrt{2}}(|\beta_{i+1}\rangle + |\beta_i\rangle) |\text{REJ}\rangle \end{aligned}$$

Now we know that  $\|\delta_i\|^2 \leq t\epsilon$  for all  $i$  (since  $|\gamma_i\rangle$  is a product state). Thus,  $\|\frac{1}{\sqrt{2}}(|\delta_{i+1}\rangle + |\delta_i\rangle)\| \leq \sqrt{2t\epsilon}$ .

Moreover,  $|\psi_+\rangle$  is a product state and so we have

$$\| |\psi_-\rangle \frac{1}{2}(|\gamma_{i+1}\rangle - |\gamma_i\rangle) + \frac{1}{\sqrt{2}}(|\delta_{i+1}\rangle + |\delta_i\rangle) \| \leq \sqrt{t\epsilon}$$

Thus,  $\| |\psi_-\rangle \frac{1}{2}(|\gamma_{i+1}\rangle - |\gamma_i\rangle) \| = \|\frac{1}{2}(|\gamma_{i+1}\rangle - |\gamma_i\rangle)\| \leq (1 + \sqrt{2})\sqrt{t\epsilon}$ .  $\square$

Then by the triangle inequality,  $\|\frac{1}{2}(|\gamma_t\rangle - |\gamma_0\rangle)\| \leq (1 + \sqrt{2})t\sqrt{t\epsilon}$ . Let  $|\Psi_\pm\rangle = \frac{1}{\sqrt{2}}(|\psi_t\rangle \pm |\psi_0\rangle)$ . The image of  $|\Psi_+\rangle = \frac{1}{\sqrt{2}}(|000\dots 0\rangle + |111\dots 1\rangle)$  is:

$$\begin{aligned} & \left( |\Psi_+\rangle \frac{1}{2}(|\gamma_t\rangle + |\gamma_0\rangle) + |\Psi_-\rangle \frac{1}{2}(|\gamma_t\rangle - |\gamma_0\rangle) + \frac{1}{\sqrt{2}}(|\delta_t\rangle + |\delta_0\rangle) \right) |\text{ACC}\rangle \\ & + \frac{1}{\sqrt{2}}(|\beta_t\rangle + |\beta_0\rangle) |\text{REJ}\rangle \end{aligned}$$

Now the trace of this state with  $P_0^{|\Psi_+\rangle}$  is the square of

$$\begin{aligned} \| |\Psi_-\rangle \frac{1}{2}(|\gamma_t\rangle - |\gamma_0\rangle) + \frac{1}{\sqrt{2}}(|\delta_t\rangle + |\delta_0\rangle) \| & \leq \| |\Psi_-\rangle \frac{1}{2}(|\gamma_t\rangle - |\gamma_0\rangle) \| + \| \frac{1}{\sqrt{2}}(|\delta_t\rangle + |\delta_0\rangle) \| \\ & \leq (1 + \sqrt{2})t\sqrt{t\epsilon} + \sqrt{2t\epsilon} \\ & \leq \sqrt{10t^3\epsilon}, \end{aligned}$$

where in the last line, we have assumed  $t \geq 2$ . That is, the iterated scheme is  $10t^3\epsilon$ -secure for  $|\Psi_+\rangle$ .  $\square$

Putting the various lemmas together, we find that, given two states  $|0\rangle$  and  $|1\rangle$  which are slightly distinguishable by the adversary, so  $D(\rho_0, \rho_1) \geq \delta$ , then in the iterated scheme,  $|000\dots 0\rangle$  and  $|111\dots 1\rangle$  are more distinguishable:  $D(\rho_{|000\dots 0\rangle}, \rho_{|111\dots 1\rangle}) \geq 1 - \eta$ , where  $\eta \leq 2 \exp(-t\delta^2/2)$ . Since the iterated scheme is  $10t^3\epsilon$ -secure for the state  $|\psi\rangle = \frac{1}{\sqrt{2}}(|000\dots 0\rangle + |111\dots 1\rangle)$ , then by the first lemma,

$$10t^3\epsilon > 1 - 2\eta \geq 1 - 4 \exp(-t\delta^2/2)$$

Choosing  $t = 1/\sqrt[3]{20\epsilon}$ , we get  $\delta \leq 4\epsilon^{1/6}$ .

**Corollary 22** *A QAS with error  $\epsilon$  requires at least  $2m(1 - \text{poly}(\epsilon))$  classical key bits.*

*Proof (of Corollary 6):* The argument is similar to the argument that  $2m$  bits of key are required for perfect encryption. We show that transmitting the key through a channel allows the transmission of almost  $2m$  bits of information.

We can consider four subsystems, two held by Alice and two held by Bob. Bob holds both halves of  $m$  Bell states (the subsystems  $B_1$  and  $B_2$ ), except that  $B_1$  has been encrypted by a key  $k$  (subsystem  $K$ ) held by Alice. Alice also holds  $R$ , a purification of the other three systems.

Using superdense coding, Bob's two systems  $B_1$  and  $B_2$  can encode  $2m$  classical bits of information. In order to recover that information, Bob needs Alice's key (system  $K$ ). Since the encryption is not perfect, however, Bob may have a small amount of information about the encoded state.

Let us imagine that Bob's systems initially encode the classical message  $000\dots 0$ . Suppose Alice wishes to send Bob the message  $M$ . Since the encryption is almost perfect, Bob's two density matrices  $\rho_B(000\dots 0)$  and  $\rho_B(M)$  are almost indistinguishable. Therefore, by the argument proving bit commitment is impossible,

Alice can change the pure state corresponding to encrypted  $000\dots 0$  to something very close to the pure state corresponding to encrypted  $M$ .

If Alice now sends  $K$  to Bob, he is able to (almost always) decode the message  $M$ . His failure probability is a polynomial in  $\epsilon$ , so he has received  $2m(1 - \text{poly}(\epsilon))$  bits of information, and therefore  $K$  must consist of at least  $2m(1 - \text{poly}(\epsilon))$  classical bits or half as many qubits.

In fact,  $K$  might as well be classical: Bob's decoding method will be to immediately measure  $K$ , since he is expecting a classical key, and therefore Alice might as well have measured  $K$  before sending it; naturally, this actually means she includes entangled qubits in the purification  $R$ . We thus restrict  $K$  to classical bits and prove the corollary.  $\square$