

# Quantum non-malleability and authentication

Gorjan Alagic and Christian Majenz

QMATH, Department of Mathematical Sciences  
University of Copenhagen

galagic@gmail.com      majenz@math.ku.dk

**Abstract.** In encryption, non-malleability is a highly desirable property: it ensures that adversaries cannot manipulate the plaintext by acting on the ciphertext. In [6], Ambainis et al. gave a definition of non-malleability for the encryption of quantum data. In this work, we show that this definition is too weak, as it allows adversaries to “inject” plaintexts of their choice into the ciphertext. We give a new definition of quantum non-malleability which resolves this problem. Our definition is expressed in terms of entropic quantities, considers stronger adversaries, and does not assume secrecy. Rather, we prove that *quantum non-malleability implies secrecy*; this is in stark contrast to the classical setting, where the two properties are completely independent. For unitary schemes, our notion of non-malleability is equivalent to encryption with a two-design (and hence also to the definition of [6]).

Our techniques also yield new results regarding the closely-related task of quantum authentication. We show that “total authentication” (a notion recently proposed by Garg et al. [18]) can be satisfied with two-designs, a significant improvement over the eight-design construction of [18]. We also show that, under a mild adaptation of the rejection procedure, both total authentication and our notion of non-malleability yield quantum authentication as defined by Dupuis et al. [16].

## 1 Introduction

**Background.** In its most basic form, encryption ensures secrecy in the presence of eavesdroppers. Besides secrecy, another desirable property is *non-malleability*, which guarantees that an active adversary cannot modify the plaintext by manipulating the ciphertext. In the classical setting, secrecy and non-malleability are independent: there are schemes which satisfy secrecy but are malleable, and schemes which are non-malleable but transmit the plaintext in the clear. If both secrecy and non-malleability is desired, then pairwise-independent permutations provide information-theoretically perfect (one-time) security [20]. In the computational security setting, non-malleability can be achieved by MACs, and ensures chosen-ciphertext security for authenticated encryption.

In the setting of quantum information, encryption is the task of transmitting quantum states over a completely insecure quantum channel. Information-theoretic secrecy for quantum encryption is well-understood. Non-malleability,

on the other hand, has only been studied in one previous work, by Ambainis, Bouda and Winter [6]. Their definition (which we will call ABW-non-malleability, or ABW-NM) requires that the scheme satisfies secrecy, and that the “effective channel”  $\text{Dec} \circ A \circ \text{Enc}$  of any adversary  $A$  amounts to either the identity map or replacement by some fixed state. In the case of unitary schemes, ABW-NM is equivalent to encrypting with a unitary two-design. Unitary two-designs are a natural quantum analogue of pairwise-independent permutations, and can be efficiently constructed in a number of ways (see, e.g., [10, 14].)

While quantum non-malleability has only been considered by [6], the closely-related task of quantum authentication (where decryption is allowed to reject) has received significant attention (see, e.g., [2, 7, 11, 16, 18].) The widely-adopted definition of Dupuis, Nielsen and Salvail asks that the averaged effective channel of any adversary is close to a map which does not touch the plaintext [16]; we refer to this notion as DNS-authentication. Recent work by Garg, Yuen and Zhandry [18] established another notion of quantum authentication, which they call “total authentication.” The notion of total authentication has two major differences from previous definitions: (i.) it asks for success with high probability over the choice of keys, rather than simply on average, and (ii.) it makes no demands whatsoever in the case that decryption rejects. We refer to this notion of quantum authentication as GYZ-authentication. In [18], it is shown that GYZ-authentication can be satisfied with unitary eight-designs.

**This work.** In this work, we devise a new definition of non-malleability (denoted NM) for quantum encryption, improving on ABW-NM in a number of ways. First, our definition is expressed in terms of entropic quantities, which allows us to bring several quantum-information-theoretic techniques to bear (such as decoupling.) Second, we consider more powerful adversaries, which can possess side information about the plaintext. Third, we remove the possibility of a “plaintext injection” attack, whereby an adversary against an ABW-NM scheme can send a plaintext of their choice to the receiver. Finally, our definition does not demand secrecy; instead, we show that *quantum secrecy is a consequence of quantum non-malleability*. This is a significant departure from the classical case, and is analogous to the fact that quantum authentication implies secrecy [7].

The primary consequence of our work is twofold: first, encryption with unitary two-designs satisfies all of the above notions of quantum non-malleability; second, when equipped with blank “tag” qubits, the same scheme also satisfies all of the above notions of quantum authentication. A more detailed summary of the results is as follows. For schemes which have unitary encryption maps, we prove that NM is equivalent to encryption with unitary two-designs, and hence also to ABW-NM. For non-unitary schemes, we prove a characterization theorem for NM schemes that shows that NM implies ABW-NM, and provide a strong separation example between NM and ABW-NM (the aforementioned plaintext injection attack). In the case of GYZ authentication, we prove that two-designs (with tags) are sufficient, a significant improvement over the state-of-the-art, which requires eight-designs [18]. Moreover, the simulation of adversaries in this

proof is efficient, in the sense of Broadbent and Wainwright [11]. Finally, we show that GYZauthentication implies DNS-authentication, and that equipping an arbitrary NM scheme with tags yields DNS-authentication.

We remark that, after the initial version of our results was submitted, an independent work of C. Portmann gave an alternative proof that GYZ-authentication can be satisfied by the 2-design scheme [26].

## 1.1 Summary of contributions

In the following, all schemes are symmetric-key encryption schemes for quantum data, in the information-theoretic security setting.

**Quantum non-malleability.** We begin with non-malleability, in both the perfect setting (Section 3) and the approximate setting (Section 4).

1. **New definition of non-malleability.** We give a new definition of quantum non-malleability (NM), in terms of the information gain of an adversary's *effective attack* on the plaintext. The quantum registers are: plaintext  $A$ , ciphertext  $C$ , user's reference  $R$ , and adversary's side information  $B$ .

**Definition 1.1 (NM, informal)** *A scheme is non-malleable (NM) if for any  $\varrho_{ABR}$  and any attack  $\Lambda_{CB \rightarrow C\tilde{B}}$ , the effective attack  $\tilde{\Lambda}_{AB \rightarrow A\tilde{B}}$  satisfies*

$$I(AR : \tilde{B})_{\tilde{\Lambda}(\varrho)} \leq I(AR : B)_{\varrho} + h(p_{=}(A, \varrho)).$$

The binary entropy term is necessary because adversaries can always simply record whether they disturbed the ciphertext (see Definition 3.4).

2. **Results on non-malleability.** Our first result is an alternative characterization of NM, in terms of the form of the effective map  $\tilde{\Lambda}$ .

**Theorem 1.2 (informal)** *A scheme is NM if and only if, for any attack  $\Lambda_{CB \rightarrow C\tilde{B}}$ , there exist maps  $\Lambda'_{B \rightarrow \tilde{B}}$ ,  $\Lambda''_{B \rightarrow \tilde{B}}$  such that the effective attack satisfies*

$$\tilde{\Lambda} = \text{id}_A \otimes \Lambda' + \frac{1}{|C|^2 - 1} (|C| \langle D_K(\mathbb{1}_C) \rangle - \text{id})_A \otimes \Lambda''.$$

The fact that NM implies ABW-NM is an immediate corollary. The new definition is strictly stronger than ABW-NM: we give a scheme which is secure under ABW-NM but insecure under NM. This scheme is in fact susceptible to a powerful attack, whereby a simple adversary can replace the output of decryption with a plaintext of the adversary's choice. On the other hand, if we restrict our attention to schemes where the encryption maps are unitary, then we are able to show the following.

**Theorem 1.3 (informal)** *Let  $\Pi$  be a scheme such that encryption  $E_k$  is unitary for all keys  $k$ . Then  $\Pi$  is NM if and only if  $\{E_k\}_k$  is a two-design.*

By the results of [6], we conclude that NM and ABW-NM are in fact equivalent for unitary schemes. Finally, we show that NM implies secrecy.

**Theorem 1.4 (informal)** *Quantum non-malleability implies secrecy.*

3. **Authentication from non-malleability.** Our final result in the setting of non-malleability shows that, by adding a “tag” space to the plaintext (as in the Clifford scheme [2]), we can turn an NM scheme into an authentication scheme as defined in [16]. More precisely, given an encryption scheme  $\Pi = \{E_k\}$ , we define  $\Pi_t^{\text{tag}}$  to be a new scheme whose encryption is  $\varrho \mapsto E_k(\varrho_A \otimes |0\rangle\langle 0|_B^{\otimes t})E_k^\dagger$ , and whose decryption rejects unless  $B$  measures to  $|0^t\rangle$ .

**Theorem 1.5 (informal)** *Let  $\Pi = \{E_k\}$  be an encryption scheme. If  $\Pi$  is NM, then  $\Pi_t^{\text{tag}}$  is  $2^{2-t}$ -DNS-authenticating.*

**Quantum authentication.** Our results on quantum authentication are summarized as follows. We note that, strictly speaking, our definitions of authentication deviate slightly from the original versions [16, 18], in that decryption outputs a reject symbol in place of the plaintext (rather than setting an auxiliary bit to “reject.”) This adaptation is convenient for reasons we will return to later.

1. **GYZ implies DNS.** First, we show that GYZ-authentication implies DNS-authentication. We remark that this is not trivial: on one hand, GYZ strengthens DNS by requiring high probability of success (rather than success on-average); on the other hand, in the reject case GYZ requires nothing while DNS makes rather stringent demands. Nonetheless, we show the following.

**Theorem 1.6 (informal)** *Let  $\Pi$  be an encryption scheme. If  $\Pi$  is  $\varepsilon$ -GYZ-authenticating, then it is also  $O(\sqrt{\varepsilon})$ -DNS-authenticating.*

2. **GYZ is achievable with 2-designs.** Next, we show that GYZ-authentication is achieved with a “tagged” two-design scheme. The analysis of [18] required eight-designs for the same construction.

**Theorem 1.7 (informal)** *Let  $\Pi = \{E_k\}_k$  be a  $2^{-t}$ -approximate 2-design scheme. Then  $\Pi_t^{\text{tag}}$  is  $2^{-\Omega(t)}$ -GYZ-authenticating.*

3. **GYZ authentication from non-malleability.** As a straightforward consequence of Theorem 1.3 and Theorem 1.7, we finally record that tagging a unitary non-malleable scheme results in a GYZ-authenticating scheme.

**Corollary 1.8 (informal)** *There exists a constant  $r > 0$  such that the following holds. If  $\Pi$  is a unitary  $\Omega(2^{-rn})$ -NM scheme for  $n$ -qubit messages, and  $t = \text{poly}(n)$ , then  $\Pi_t^{\text{tag}}$  is  $2^{-\Omega(\text{poly}(n))}$ -GYZ-authenticating.*

A sufficiently strong NM scheme can be constructed via the  $\varepsilon$ -approximate version of Theorem 1.3 (see Theorem 4.5 and Remark 2.3 below.)

The remainder of the paper is structured as follows. In Section 2, we review some basic facts regarding quantum states, registers, and channels, and recall several useful facts about unitary designs. In Section 3, we consider the

exact setting, beginning with perfect secrecy and then continuing to perfect non-malleability (NM) and the relevant new results; we also discuss the relationship to ABW-NM in detail. We continue in [Section 4](#) with the approximate setting, again beginning with secrecy and then continuing to approximate non-malleability. We end with the new results on quantum authentication, in [Section 4.2](#).

## 2 Preliminaries

### 2.1 Quantum states, registers, and channels.

We assume basic familiarity with the formalism of quantum states, operators, and channels. We denote quantum registers (i.e., systems and their subsystems) with capital Latin letters, e.g.,  $A, B, C$ . The Hilbert space corresponding to system  $A$  is denoted by  $\mathcal{H}_A$ . For a register  $A$ , we denote the dimension of  $\mathcal{H}_A$  by  $|A|$ . We emphasize that, in this work, all Hilbert spaces will be finite-dimensional.

The space operators on  $\mathcal{H}_A$  is denoted  $\mathcal{B}(\mathcal{H}_A)$ . We say that a quantum state is classical if it is diagonal in the standard (i.e., computational) basis. We denote the adjoint of an operator  $X \in \mathcal{B}(\mathcal{H})$  by  $X^\dagger$  and its transpose with respect to the computational basis by  $X^T$ . Where necessary, we will write a quantum state  $\varrho \in \mathcal{B}(\mathcal{H}_A \otimes \mathcal{H}_B \otimes \mathcal{H}_C)$  as  $\varrho_{ABC}$  to emphasize that the state is a multipartite state over registers  $A, B$ , and  $C$ . When such a state has already been defined, we will write reduced states by omitting the traced-out registers, e.g.,  $\varrho_A := \text{Tr}_{BC}[\varrho_{ABC}]$ . We single out some special states which will appear frequently. Fix two systems  $S, S'$  with  $|S| = |S'|$ . We let

$$|\phi^+\rangle_{SS'} = |S|^{-1/2} \sum_i |ii\rangle_{SS'} \quad \text{and} \quad \phi_{SS'}^+ = |\phi^+\rangle\langle\phi^+|_{SS'}$$

denote the maximally entangled state on the bipartite system  $SS'$  (expressed as a pure state on the left, and as a density operator on the right.) Furthermore, we let  $\Pi_{SS'}^- = \mathbb{1}_{SS'} - \phi_{SS'}^+$  and  $\tau_{SS'}^- = \Pi_{SS'}^- / (|S|^2 - 1)$ . We also set  $\tau_S = \mathbb{1}_S / |S|$  to be the maximally mixed state on  $S$ .

We denote the von Neumann entropy of a state  $\varrho_A$  by  $H(A)_\varrho$ , and the joint entropy of  $\varrho_{AB}$  by  $H(AB)_\varrho$ . We recall that the quantum mutual information of  $\varrho_{AB}$  is defined by

$$I(A : B)_\varrho := H(A)_\varrho + H(B)_\varrho - H(AB)_\varrho.$$

The quantum conditional mutual information of  $\varrho_{ABC}$  is defined by

$$I(A : B|C)_\varrho := H(AC)_\varrho + H(BC)_\varrho - H(ABC)_\varrho - H(C)_\varrho.$$

These quantities are nonnegative [\[21\]](#) and satisfy a chain rule:

$$I(A : BC|D)_\varrho = I(A : B|D)_\varrho + I(A : C|BD)_\varrho.$$

We remark that the above also holds for trivial  $D$ . Together with the Stinespring dilation theorem [27], non-negativity [22] and the chain rule imply the data processing inequality

$$I(A : \tilde{B}|C)_{\Lambda(\varrho)} \leq I(A : B|C)_\varrho,$$

when  $\Lambda$  is a CPTP (completely-positive, trace-preserving) map from  $\mathcal{B}(\mathcal{H}_B)$  to  $\mathcal{B}(\mathcal{H}_{\tilde{B}})$ . An important special case is where  $B = B_1 B_2$  and  $\Lambda = \text{Tr}_{B_2}$  discards the contents of  $B_2$ .

We will refer to valid transformations between quantum states as channels, or CPTP maps. We will sometimes also consider trace-non-increasing completely-positive (CP) maps. When necessary, we will emphasize the input and output spaces of a map  $\Lambda : \mathcal{B}(\mathcal{H}_A \otimes \mathcal{H}_B) \rightarrow \mathcal{B}(\mathcal{H}_C)$  by writing  $\Lambda_{AB \rightarrow C}$ . We denote the identity channel on, e.g., register  $A$  by  $\text{id}_{A \rightarrow A}$  (or simply  $\text{id}_A$ ) and the channel from register  $A$  to  $A'$  with constant output  $\sigma_{A'}$  by  $\langle \sigma \rangle_{A \rightarrow A'}$ . When composing operators on many registers, and if the context allows, we will elide tensor products with the identity operator. So, for example, with  $\Lambda$  as above we may write  $\tau_{CD} = \Lambda \varrho_{ABD}$  in place of  $\tau_{CD} = (\Lambda \otimes \text{id}_D) \varrho_{ABD}$ .

A standard tool in this setting is the Choi-Jamiolkowski (CJ) isomorphism [12, 19]. Let  $\Xi_{A \rightarrow B} : \mathcal{B}(\mathcal{H}_A) \rightarrow \mathcal{B}(\mathcal{H}_B)$  be a linear operator. Then its CJ matrix is defined as

$$(\eta_\Xi)_{BA'} = \Lambda_{A \rightarrow B}(\phi_{AA'}^+). \quad (2.1)$$

The linear operator mapping  $\Xi$  to  $\eta_\Xi$  is an isomorphism of vector spaces and  $\eta_\Xi$  is positive semidefinite iff  $\Xi$  is CP. Moreover  $\Xi_{A \rightarrow B}$  is TP iff  $(\eta_\Xi)_{A'} = \tau_A$ . The inverse of the CJ isomorphism is given by the equation

$$\Xi_{A \rightarrow B}(X_A) = |A\rangle \langle A| \text{Tr}_{A'} [X_{A'}^T (\eta_\Xi)_{BA'}]. \quad (2.2)$$

We denote the swap operator by  $F : |i\rangle \otimes |j\rangle \mapsto |j\rangle \otimes |i\rangle$ .

**Lemma 2.1 (Swap trick [17])** *For matrices  $A$  and  $B$ ,  $\text{Tr}[AB] = \text{Tr}[FA \otimes B]$ .*

We will make frequent use of the trace norm  $\|\cdot\|_1$ , the operator norm  $\|\cdot\|_\infty$ , and the diamond norm  $\|\Lambda_{A \rightarrow B}\|_\diamond := \max_{\varrho_{AA'}} \|\Lambda_{A \rightarrow B} \otimes \text{id}_{A'}(\varrho_{AA'})\|_1$ ; here the max is taken over all pure quantum states  $\varrho_{AA'}$  and  $\mathcal{H}_A \cong \mathcal{H}_{A'}$ . Recall that the Hölder inequality for operators states that, for any two operators  $X$  and  $Y$ ,

$$\text{Tr}[XY] \leq \|XY\|_1 \leq \|X\|_1 \|Y\|_\infty. \quad (2.3)$$

## 2.2 Unitary designs.

We now recall the definition of unitary  $t$ -design, and some relevant variants. We begin by considering three different types of “twirls.”

1. For a finite subset  $D \subset U(\mathcal{H})$  of the unitary group on some finite dimensional Hilbert space  $\mathcal{H}$ , let

$$\mathcal{T}_D^{(t)}(X) = \frac{1}{|D|} \sum_{U \in D} U^{\otimes t} X (U^\dagger)^{\otimes t} \quad (2.4)$$

be the associated  $t$ -twirling channel. If we take the entire unitary group (rather than just a finite subset), then we get the Haar  $t$ -twirling channel

$$\mathcal{T}_{\text{Haar}}^{(t)}(X) = \int U^{\otimes t} X (U^\dagger)^{\otimes t} dU. \quad (2.5)$$

2. We define the  $U$ - $\bar{U}$  twirl with respect to finite  $D \subset U(\mathcal{H})$  by

$$\bar{\mathcal{T}}_D(X) = \frac{1}{|D|} \sum_{U \in D} (U \otimes \bar{U}) X (U \otimes \bar{U})^\dagger. \quad (2.6)$$

The analogous  $U$ - $\bar{U}$  Haar twirling channel is denoted by  $\bar{\mathcal{T}}_{\text{Haar}}$ .

3. The third notion is called a channel twirl, and is defined in terms of  $U$ - $\bar{U}$ -twirling. Given a channel  $\Lambda$ , let  $\eta_\Lambda$  be the CJ state of  $\Lambda$ . The channel twirl  $\bar{\mathcal{T}}_D^{ch}(\Lambda)$  of  $\Lambda$  is defined to be the channel whose CJ state is  $\bar{\mathcal{T}}_D(\eta_\Lambda)$ .

Next, we define the three corresponding notions of designs.

**Definition 2.2** *Let  $D \subset U(\mathcal{H})$  be a finite set. We define the following.*

- If  $\|\mathcal{T}_D^{(t)} - \mathcal{T}_{\text{Haar}}^{(t)}\|_\diamond \leq \delta$  holds, then  $D$  is a  $\delta$ -approximate  $t$ -design.
- If  $\|\bar{\mathcal{T}}_D - \bar{\mathcal{T}}_{\text{Haar}}\|_\diamond \leq \delta$  holds, then  $D$  is a  $\delta$ -approximate  $U$ - $\bar{U}$ -twirl design.
- If  $\|\bar{\mathcal{T}}_D^{ch}(\Lambda) - \bar{\mathcal{T}}_{\text{Haar}}^{ch}(\Lambda)\|_\diamond \leq \delta$  holds for all CPTP maps  $\Lambda$ , then  $D$  is a  $\delta$ -approximate channel-twirl design.

For all three of the above, the case  $\delta = 0$  is called an “exact design” (or simply “design”.) All three notions of design are equivalent in the exact case. In the approximate case they are still connected, but there are some nontrivial costs in the approximation quality (See [23], Lemma 2.2.14, and an additional easy lemma proven in the full version [3]).

It is well-known that  $\varepsilon$ -approximate  $t$ -designs on  $n$  qubits can be generated by random quantum circuits of size polynomial in  $n, t$  and  $\log(1/\varepsilon)$  [10]. In particular, the size of these circuits is polynomial even for exponentially-small choices of  $\varepsilon$ . We emphasize this observation as follows.

**Remark 2.3** *Fix a polynomial  $t$  in  $n$ . Then, for any  $\varepsilon > 0$ , a random  $n$ -qubit quantum circuit consisting of  $\text{poly}(n, \log(1/\varepsilon))$  gates (from a universal set) satisfies every notion of  $\varepsilon$ -approximate  $t$ -design in Definition 2.2.*

For exact designs, we point out two important constructions. First, the prototypical example of a unitary one-design on  $n$  qubits is the  $n$ -qubit Pauli group. For exact unitary two-designs, the standard example is the Clifford group, which is the normalizer of the  $n$ -qubit Pauli group. Alternatively, the Clifford group is generated by circuits from the gate set  $\{H, P, \text{CNOT}\}$ . It is well-known that one can efficiently generate exact unitary two-designs on  $n$ -qubits by building appropriate circuits from this gate set, using  $O(n^2)$  random bits [1, 14].

### 3 The zero-error setting

We begin with the zero-error. In the case of secrecy, zero-error means that schemes cannot leak any information whatsoever. In the case of non-malleability, zero-error means that the adversary cannot increase their correlations with the secret by even an infinitesimal amount (except by trivial means; see below.)

#### 3.1 Perfect secrecy

We begin with a definition of symmetric-key quantum encryption. Our formulation treats rejection during decryption in a slightly different manner from previous literature.

**Definition 3.1 (Encryption scheme)** *A symmetric-key quantum encryption scheme (QES) is a triple  $(\tau_K, E, D)$  consisting of a classical state  $\tau_K \in \mathcal{B}(\mathcal{H}_K)$  and a pair of channels*

$$\begin{aligned} E &: \mathcal{B}(\mathcal{H}_A \otimes \mathcal{H}_K) \longrightarrow \mathcal{B}(\mathcal{H}_C \otimes \mathcal{H}_K) \\ D &: \mathcal{B}(\mathcal{H}_C \otimes \mathcal{H}_K) \longrightarrow \mathcal{B}((\mathcal{H}_A \oplus \mathbb{C}|\perp\rangle) \otimes \mathcal{H}_K) \end{aligned}$$

*satisfying  $[D \circ E](\cdot \otimes |k\rangle\langle k|) = (\text{id}_A \oplus 0_\perp) \otimes |k\rangle\langle k|$  for all  $k$ .*

The Hilbert spaces  $\mathcal{H}_A$ ,  $\mathcal{H}_C$  and  $\mathcal{H}_K$  are implicitly given by the triple  $(\tau_K, E, D)$ . The state  $|\perp\rangle$  is an error flag that allows the decryption map to report an error. For notational convenience when dealing with these schemes, we set

$$\begin{aligned} E_k &= E(\cdot \otimes |k\rangle\langle k|) & E_K &= \text{Tr}_K E(\cdot \otimes \tau_K) \\ D_k &= D(\cdot \otimes |k\rangle\langle k|) & D_K &= \text{Tr}_K D(\cdot \otimes \tau_K). \end{aligned}$$

We will often slightly abuse notation by referring to decryption maps  $D_k$  as maps from  $C$  to  $A$ ; in fact, the output space of  $D_k$  is really the slightly larger space  $\bar{A} := A \oplus \mathbb{C}|\perp\rangle$ .

It is natural to define secrecy in the quantum world in terms of quantum mutual information. However, instead of asking for the ciphertext to be uncorrelated with the plaintext as in the classical case, we ask for the ciphertext to be uncorrelated from any reference system.

**Definition 3.2 (Perfect secrecy)** *A QES  $(\tau_K, E, D)$  satisfies information - theoretic secrecy (ITS) if, for any Hilbert space  $\mathcal{H}_B$  and any  $\varrho_{AB} \in \mathcal{B}(\mathcal{H}_A \otimes \mathcal{H}_B)$ , setting  $\sigma_{CBK} = E(\varrho_{AB} \otimes \tau_K)$  implies  $I(C : B)_\sigma = 0$ .*

We note that, for perfect ITS, adding side information is unnecessary: the definition already implies that the ciphertext is in product with *any* other system. In particular, if the adversary has some auxiliary system  $E$  in their possession, then  $I(B : CE)_\sigma = I(B : E)_\sigma$ . Several definitions of secrecy for symmetric-key quantum encryption have appeared in the literature, but the above formulation



appears to be new. It can be shown that ITS is equivalent to perfect indistinguishability of ciphertexts (IND). The latter notion is a special case of an early indistinguishability-based definition of Ambainis et al. [5].

In many situations it makes sense to restrict ourselves to QES that have identical plaintext and ciphertext spaces; due to correctness, this is equivalent to unitarity.

**Definition 3.3 (Unitary scheme)** A QES  $(\tau_K, E, D)$  is called unitary if the encryption and decryption maps are controlled unitaries, i.e., if there exists  $V = \sum_k U_A^{(k)} \otimes |k\rangle\langle k|_K$  such that  $E(X) = VXV^\dagger$ .

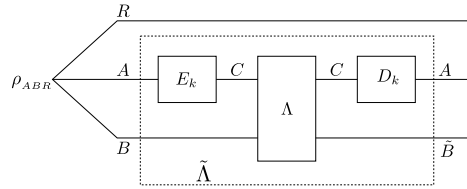
It is straightforward to prove that, for unitary schemes, ITS is equivalent to the statement that the encryption maps  $\{E_k\}$  form a unitary 1-design. Note that unitarity of  $E_k$  and correctness imply unitarity of  $D_k$ .

### 3.2 A new notion of non-malleability

**Definition.** We consider a scenario involving a user Alice and an adversary Mallory. The scenario begins with Mallory preparing a tripartite state  $\rho_{ABR}$  over three registers: the plaintext  $A$ , the reference  $R$ , and the side-information  $B$ . The registers  $A$  and  $R$  are given to Alice, while Mallory keeps  $B$ . Alice then encrypts  $A$  into a ciphertext  $C$  and then transmits (or stores) it in the open. Mallory now applies an attack map

$$\Lambda : \mathcal{B}(\mathcal{H}_C \otimes \mathcal{H}_B) \rightarrow \mathcal{B}(\mathcal{H}_C \otimes \mathcal{H}_{\tilde{B}}).$$

Mallory keeps the (transformed) side-information  $\tilde{B}$  and returns  $C$  to Alice. Finally, Alice decrypts  $C$  back to  $A$ , and the scenario ends. We are now interested



**Fig. 1.** The quantum non-malleability scenario.

in measuring the extent to which Mallory was able to increase her correlations with Alice's systems  $A$  and  $R$ . This can be understood by analyzing the mutual information  $I(AR : \tilde{B})_{\tilde{\Lambda}(\rho)}$  where  $\tilde{\Lambda}_{AB \rightarrow A\tilde{B}}$  is the *effective channel* corresponding to Mallory's attack:

$$\tilde{\Lambda} = \text{Tr}_K(D \circ \Lambda \circ E)((\cdot) \otimes \tau_K). \quad (3.1)$$

We point out one way in which Mallory can always increase these correlations, regardless of the structure of the encryption scheme. First, she flips a coin  $b$ , and records the outcome in  $B$ . If  $b = 1$ , she replaces the contents of  $C$  with some fixed state  $\sigma_C$ , and otherwise she leaves  $C$  untouched. One then sees that Mallory's correlations have increased by  $h(p_=(A, \varrho))$ , where  $h$  denotes binary entropy and  $p_=(A, \varrho)$  is defined as follows.

$$p_=(A, \varrho) = \text{Tr} [(\phi_{CC'}^+ \otimes \mathbb{1}_{\tilde{B}})A(\phi_{CC'}^+ \otimes \varrho_B)] . \quad (3.2)$$

This quantity is the inner product between the identity map and the map  $A((\cdot) \otimes \varrho_B)$ , expressed in terms of CJ states. Intuitively, it measures the probability with which Mallory chooses to apply the identity map; taking the binary entropy then gives us the information gain resulting from recording this choice.

We are now ready to define information-theoretic quantum non-malleability. Stated informally, a scheme is non-malleable if Mallory can only implement the attacks described above.

**Definition 3.4 (Non-malleability)** *A QES  $(\tau_K, E, D)$  is non-malleable (NM) if for any state  $\varrho_{ABR}$  and any CPTP map  $\Lambda_{CB \rightarrow C\tilde{B}}$ , we have*

$$I(AR : \tilde{B})_{\tilde{\Lambda}(\varrho)} \leq I(AR : B)_\varrho + h(p_=(A, \varrho)) . \quad (3.3)$$

One might justifiably wonder if the term  $h(p_=(A, \varrho))$  is too generous to the adversary. However, as we showed above, every scheme is vulnerable to an attack which gains this amount of information. This term also appears (somewhat disguised) in the classical setting. In fact, if a classical encryption scheme satisfies [Definition 3.4](#) against classical adversaries, then it also satisfies classical information-theoretic non-malleability as defined in [\[20\]](#).

[Definition 3.4](#) directly generalizes the classical information-theoretic definition from [\[20\]](#). In some settings, it might be preferable to have a definition which characterizes the set of effective attack channels as was done in [\[6\]](#). As it turns out, NM can be defined in this way.

**Theorem 3.7 (Non-malleability, alternative form)** *A QES  $(\tau, E, D)$  is NM if and only if for any attack  $\Lambda_{CB \rightarrow C\tilde{B}}$ , the effective map  $\tilde{\Lambda}_{AB \rightarrow A\tilde{B}}$  has the form*

$$\tilde{\Lambda} = \text{id}_A \otimes A'_{B \rightarrow \tilde{B}} + \frac{1}{|C|^2 - 1} (|C|^2 \langle D_K(\tau) \rangle - \text{id})_A \otimes A''_{B \rightarrow \tilde{B}} \quad (3.4)$$

where  $A' = \text{Tr}_{CC'}[\phi_{CC'}^+ A(\phi_{CC'}^+ \otimes (\cdot))] and  $A'' = \text{Tr}_{CC'}[\Pi_{CC'}^- A(\phi_{CC'}^+ \otimes (\cdot))]$ .$

The proof of this theorem is postponed to the results section below (proof sketch) and the appendix.

Finally, as we will show in later sections, [Definition 3.4](#) implies ABW-NM (see [Definition 3.8](#)), and schemes satisfying [Definition 3.4](#) are sufficient for building quantum authentication under the strongest known definitions.

**Non-malleability implies secrecy.** In the classical case, non-malleability is independent from secrecy: the one-time pad is secret but malleable, and non-malleability is unaffected by appending the plaintext to each ciphertext. In the quantum case, on the other hand, we can show that NM implies secrecy. This is analogous to the fact that “quantum authentication implies encryption” [7].

**Proposition 3.5** *Let  $(\tau_K, E, D)$  be an NM QES. Then  $(\tau_K, E, D)$  is ITS.*

*Proof.* Let  $B, \varrho_{AB}$ , and  $\sigma_{CBK} = E(\varrho_{AB} \otimes \tau_K)$  be as in the definition of ITS (Definition 3.2). We first rename  $B$  to  $R$ . We then consider the non-malleability property in the following special-case scenario. The initial side-information register is empty, the final side-information register  $\tilde{B}$  satisfies  $\mathcal{H}_{\tilde{B}} \cong \mathcal{H}_C$ , and the adversary map  $A_{C \rightarrow C\tilde{B}}$  is defined as follows. Note that the “ciphertext-extraction” map  $\Theta_{C \rightarrow C\tilde{B}} = \text{id}_{C \rightarrow \tilde{B}}(\cdot) \otimes \tau_C$  has CJ state  $\eta_{CC'\tilde{B}}^\Theta = \phi_{C'\tilde{B}}^+ \otimes \tau_C$ . We choose  $\Lambda$  so that its CJ state satisfies

$$\eta_{CC'\tilde{B}}^\Lambda = \frac{d^2}{d^2 - 1} \Pi_{CC'}^- \eta_{CC'\tilde{B}}^\Theta \Pi_{CC'}^-. \quad (3.5)$$

Applying the above projection to the CJ state of  $\Theta$  ensures that  $\Lambda$  will have  $p_=(\Lambda) = 0$  (note:  $p_=(\Theta) > 0$ .)

Direct calculation of the  $C'\tilde{B}$  marginal of the CJ state of  $\Lambda$  yields

$$\eta_{C'\tilde{B}}^\Lambda = \frac{d^2 - 2}{d^2 - 1} \phi_{C'\tilde{B}}^+ + \frac{1}{d^2 - 1} \tau_{C'} \otimes \tau_{\tilde{B}}. \quad (3.6)$$

This implies that the output  $\sigma_{AR\tilde{B}} = \tilde{\Lambda}_{A \rightarrow A\tilde{B}}(\varrho_{AB})$  of the effective channel  $\tilde{\Lambda}$  will satisfy

$$\sigma_{\tilde{B}R} = \frac{d^2 - 2}{d^2 - 1} \gamma_{\tilde{B}R} + \frac{1}{d^2 - 1} \tau_{\tilde{B}} \otimes \varrho_R, \quad (3.7)$$

where  $\gamma_{CR} = (E_K)_{A \rightarrow C}(\varrho_{AR})$  and we used the fact that  $\mathcal{H}_{\tilde{B}} \cong \mathcal{H}_C$ . By non-malleability, we have

$$I(\tilde{B} : R)_\sigma + I(\tilde{B} : A|R)_\sigma = I(\tilde{B} : AR)_\sigma = 0. \quad (3.8)$$

In particular,  $I(\tilde{B} : R)_\sigma = 0$  and thus  $\sigma_{\tilde{B}R} = \sigma_{\tilde{B}} \otimes \varrho_R$ . It follows by Equation (3.7) that

$$\gamma_{\tilde{B}R} = \frac{d^2 - 1}{d^2 - 2} \left( \sigma_{\tilde{B}} - \frac{1}{d^2 - 1} \tau_{\tilde{B}} \right) \otimes \varrho_R, \quad (3.9)$$

i.e.,  $\gamma_{\tilde{B}R}$  is a product state. This is precisely the definition of information-theoretic secrecy.  $\square$

**Characterization of non-malleable schemes.** Next, we provide a characterization of non-malleable schemes. First, we show that unitary schemes are equivalent to encryption with a unitary 2-design.

**Theorem 3.6** *A unitary QES  $(\tau_K, E, D)$  is NM if and only if  $\{E_k\}_{k \in K}$  is a unitary 2-design.*

This fact is particularly intuitive when the 2-design is the Clifford group, a well-known exact 2-design. In that case, a Pauli operator acting on only one ciphertext qubit will be “propagated” (by the encryption circuit) to a completely random Pauli on all plaintext qubits. The plaintext is then maximally mixed, and the adversary gains no information. The Clifford group thus yields a perfectly non-malleable (and perfectly secret) encryption scheme using  $O(n^2)$  bits of key [1].

It will be convenient to prove [Theorem 3.6](#) as a consequence of our general characterization theorem, which is as follows.

**Theorem 3.7** *Let  $(\tau, E, D)$  be a QES. Then  $(\tau, E, D)$  is NM if and only if, for any attack  $\Lambda_{CB \rightarrow C\bar{B}}$ , the effective map  $\tilde{\Lambda}_{AB \rightarrow A\bar{B}}$  has the form*

$$\tilde{\Lambda} = \text{id}_A \otimes A'_{B \rightarrow \bar{B}} + \frac{1}{|C|^2 - 1} (|C|^2 \langle D_K(\tau) \rangle - \text{id})_A \otimes A''_{B \rightarrow \bar{B}} \quad (3.10)$$

where  $A' = \text{Tr}_{CC'}[\phi_{CC'}^+ \Lambda(\phi_{CC'}^+ \otimes (\cdot))]$  and  $A'' = \text{Tr}_{CC'}[\Pi_{CC'}^- \Lambda(\phi_{CC'}^+ \otimes (\cdot))]$ .

We remark that the forward direction holds even if  $(\tau, E, D)$  only fulfills the NM condition (Equation (3.3)) against adversaries with empty side-information  $B$ . The proof of [Theorem 3.7](#) (with this strengthening) is sketched below. The full proof is somewhat technical and can be found in [Appendix B](#). More precisely, we prove the stronger [Theorem B.3](#), which implies the above by setting  $\varepsilon = 0$ .

*Proof sketch.* The first implication, i.e. NM implies Equation (3.10), is best proven in the Choi-Jamioilkowski picture. Here, any QES defines a map

$$\mathcal{E}_{CC' \rightarrow AA'} = \frac{1}{|K|} \sum_k D_k \otimes E_k^T, \quad (3.11)$$

where the transpose  $E_k^T$  is the map whose Kraus operators are the transposes of the Kraus operators of  $E_k$  (in the standard basis). Our goal is to prove that this map essentially acts like the  $U\bar{U}$ -twirl. We decompose the space  $\mathcal{H}_C^{\otimes 2}$  as

$$\mathcal{H}_C^{\otimes 2} = \mathbb{C}|\phi^+\rangle \oplus \text{supp}\Pi^- \quad (3.12)$$

which induces a decomposition of

$$\begin{aligned} \mathcal{B}(\mathcal{H}_C^{\otimes 2}) &= \mathbb{C}|\phi^+\rangle\langle\phi^+| \oplus \left\{ |\phi^+\rangle\langle v| \mid \langle\phi^+ | v\rangle = 0 \right\} \\ &\oplus \left\{ |v\rangle\langle\phi^+| \mid \langle\phi^+ | v\rangle = 0 \right\} \oplus \left\{ X \in B \mid \langle\phi^+ | X = X|\phi^+\rangle = 0 \right\}. \end{aligned} \quad (3.13)$$

On the first and last direct summands, the correct behavior of  $\mathcal{E}$  is easy to show: the first one corresponds to the identity, and the last one to the non-identity channels  $\Lambda$  with  $p_-(\Lambda) = 0$ . For the remaining two spaces, we employ [Lemma A.3](#) which shows that the encryption map of any valid encryption scheme has the form of appending an ancillary mixed state and then applying an isometry. Evaluating  $\mathcal{E}(|\phi^+\rangle\langle v|)$  for  $\langle\phi^+ | v\rangle = 0$  reduces to evaluating the adjoint of the average encryption map,  $E_K^\dagger$ , on traceless matrices. It is, however, easy to verify that

$$\text{Tr}_A \mathcal{E}_{CC' \rightarrow AA'}(\sigma_C \otimes (\cdot)_{C'}) = (E_K^T)_{C' \rightarrow A'}$$

for any  $\sigma_C$ . This can be used to prove  $E_K = \langle \tau_C \rangle$  by observing that  $\langle \phi^+ |_{CC'} \sigma_C \otimes \varrho_{C'} | \phi^+ \rangle_{CC'} = \text{Tr}(\sigma_C \varrho_C)$ , so for rank-deficient  $\varrho$  we can calculate  $\mathcal{E}_{CC' \rightarrow AA'}(\sigma_C \otimes (\cdot)_{C'})$  using what we have already proven.

The other direction is proven by a simple application of [Lemma A.2](#).  $\square$

The fact that NM is equivalent to 2-designs (for unitary schemes) is a straightforward consequence of the above.

*Proof.* (of [Theorem 3.6](#)) First, assume  $(\tau_K, E, D)$  is a unitary NM QES with  $E_k = U_k(\cdot)U_k^\dagger$ . Then it has  $|C| = |A|$ , and  $D_K(\tau_C) = \tau_A$ , so the conclusion of [Theorem 3.7](#) in this case (i.e., Equation (3.10)) is exactly the condition for  $\{U_k\}$  to be an exact channel twirl design and therefore an exact 2-design. If  $(\tau_K, E, D)$ , on the other hand, is a unitary QES and  $\{U_k\}$  is a 2-design, then Equation (3.10) holds and the scheme is therefore NM according to [Theorem 3.7](#).

**Relationship to ABW non-malleability.** Ambainis, Bouda and Winter give a different definition of non-malleability, expressed in terms of the effective maps that an adversary can apply to the plaintext by acting on the ciphertext produced from encrypting with a random key [6]. According to their definition, a scheme is non-malleable if the adversary can only apply maps from a very restricted class *when averaging over the key, and without giving side information to the active adversary*. Let us recall their definition here.

First, given a QES  $(\tau_K, E, D)$ , we define the set  $S := \{D_K(\sigma_C) \mid \sigma_C \in \mathcal{B}(\mathcal{H}_C)\}$  consisting of all valid average decryptions. We then define the class  $C_A^S$  of all “replacement channels”. This is the set of CPTP maps belonging to the space

$$\text{span}_{\mathbb{R}}\{\text{id}_A, (X \mapsto \text{Tr}(X)\sigma_A) : \sigma_A \in S\}. \quad (3.14)$$

We then make the following definition, which first appeared in [6].

**Definition 3.8 (ABW non-malleability)** *A QES  $(\tau_K, E, D)$  is ABW-non-malleable (ABW-NM) if it is ITS, and for all channels  $\Lambda_{C \rightarrow C}$ , we have*

$$\text{Tr}_K [D_{CK \rightarrow AK} \circ \Lambda_{C \rightarrow C} \circ E_{AK \rightarrow CK}(\cdot \otimes \tau_K)] \in C_A^S. \quad (3.15)$$

As indicated in [6], an approximate version of Equation (3.15) is obtained by considering the diamond-norm distance between the effective channel and the set  $C_A^S$ ; this implies the possibility of an auxiliary reference system, which is denoted  $R$  in NM. We emphasize that this reference system is not under the control of the adversary. In particular, ABW-NM does not allow for adversaries which maintain *and actively use* side information about the plaintext system.

Another notable distinction is that [6] includes a secrecy assumption in the definition of an encryption scheme; under this assumption, it is shown that a unitary QES is ABW-NM if and only if the encryption unitaries form a 2-design. By our [Theorem 3.6](#), we see that NM and ABW-NM are equivalent in the case of unitary schemes. So, in that case, ABW-NM actually ensures a much stronger security notion than originally considered by the authors of [6].

In the general case, NM is strictly stronger than ABW-NM. First, by comparing the conditions of [Definition 3.8](#) to Equation (3.10), we immediately get the following corollary of [Theorem 3.7](#).

**Corollary 3.9** *If a QES satisfies NM, then it also satisfies ABW-NM.*

Second, we give a separation example which shows that ABW-NM is highly insecure; in fact, it allows the adversary to “inject” a plaintext of their choice into the ciphertext. This is insecure even under the classical definition of information-theoretic non-malleability of [20]. We now describe the scheme and this attack.

**Example 3.10** *Suppose  $(\tau_K, E, D)$  is a QES that is both NM and ABW-NM. Define a modified scheme  $(\tau_K, E', D')$ , with enlarged ciphertext space  $\mathcal{H}_{C'} = \mathcal{H}_C \oplus \mathcal{H}_{\hat{A}}$  (where  $\mathcal{H}_{\hat{A}} \cong \mathcal{H}_A$ ) and encryption and decryption defined by*

$$\begin{aligned} E'(X) &= E(X)_C \oplus 0_{\hat{A}} \\ D'(X) &= D_{CK \rightarrow AK}(\Pi_C X \Pi_C) + \text{id}_{\hat{A}K \rightarrow AK}(\Pi_{\hat{A}} X \Pi_{\hat{A}}). \end{aligned}$$

*Then  $(\tau_K, E', D')$  is ABW-NM but not NM.*

While encryption ignores  $\mathcal{H}_{\hat{A}}$ , decryption measures if we are in  $C$  or  $\hat{A}$  and then decrypts (in the first case) or just outputs the contents (in the second case.) This is a dramatic violation of NM: set  $\mathcal{H}_{\tilde{B}} \cong \mathcal{H}_A$ , trivial  $B$  and  $R$ , and

$$\Lambda_{C' \rightarrow C' \tilde{B}}(X) = \text{Tr}(X) 0_C \oplus |\phi^+\rangle \langle \phi^+|_{\hat{A} \tilde{B}}; \quad (3.16)$$

it follows that, for all  $\varrho$ ,

$$I(AR : \tilde{B})_{\hat{A}(\varrho)} = 2 \log |A| \gg h(|C'|^{-2}) = h(p_{=}(A, \varrho)). \quad (3.17)$$

Now let us show that  $(\tau, E', D')$  is still ABW-NM. Let  $\Lambda_{C' \rightarrow C'}$  be an attack, i.e., an arbitrary CPTP map. Then the effective plaintext map is

$$\tilde{\Lambda}_{A \rightarrow A} = D \circ \Lambda_{C \rightarrow C}^C \circ E + \Lambda_{\hat{C} \rightarrow A}^{\hat{A}} \circ E, \quad (3.18)$$

where  $\Lambda^C(X_C) = \Pi_C \Lambda(X_C \oplus 0_{\hat{A}}) \Pi_C$  and  $\Lambda^{\hat{A}}(X_C) = \text{id}_{\hat{A} \rightarrow A}(\Pi_{\hat{A}} \Lambda(X_C \oplus 0_{\hat{A}}) \Pi_{\hat{A}})$ . Since  $(\tau, E, D)$  is ITS ([Theorem 3.5](#)), there exists a fixed state  $\varrho_C^0$  such that  $E_K(\varrho_A) = \varrho_C^0$  for all  $\varrho_A$ . Since  $(\tau, E, D)$  is ABW-NM, we also know that

$$\text{Tr}_K \circ D \circ \Lambda_{C \rightarrow C}^C \circ E = \tilde{\Lambda}_1 \in C_A^S,$$

with  $S = \{D_K(\sigma_C) \mid \sigma_C \in \mathcal{B}(\mathcal{H}_C)\}$ . We therefore get

$$\tilde{\Lambda}_{A \rightarrow A} = \tilde{\Lambda}_1 + \langle \Lambda^{\hat{A}}(\varrho_C^0) \rangle \in C_A^{S'}, \quad (3.19)$$

with  $S' = \{D'_K(\sigma_{C'}) \mid \sigma_{C'} \in \mathcal{B}(\mathcal{H}_{C'})\}$ . This is true because  $S'$  contains all constant maps, as  $D'_K(0_C \oplus \varrho_{\hat{A}}) = \varrho_A$ .

## 4 The approximate setting

We now consider the case of approximate non-malleability. Approximate schemes are relevant for several reasons. First, an approximate scheme with negligible error can be more efficient than an exact one: the most efficient construction of an exact 2-design requires a quantum circuit of  $O(n \log n \log \log n)$  gates [13], where approximate 2-designs can be achieved with linear-length circuits [14]. Second, in practice, absolutely perfect implementation of all quantum gates is too much to expect—even with error-correction. Third, when passing to authentication one must allow for errors, as it is always possible for the adversary to escape detection (with low probability) by guessing the secret key.

For all these reasons, it is important to understand what happens when the perfect secrecy and perfect non-malleability requirements are slightly relaxed. In this section, we show that our definitions and results are stable under such relaxations, and prove several additional results for quantum authentication. We begin with the approximate-case analogue of perfect secrecy.

**Definition 4.1 (Approximate secrecy)** Fix  $\varepsilon > 0$ . A QES  $(\tau_K, E, D)$  is  $\varepsilon$ -approximately secret ( $\varepsilon$ -ITS) if, for any  $\mathcal{H}_B$  and any  $\varrho_{AB}$ , setting  $\sigma_{CBK} = E(\varrho_{AB} \otimes \tau_K)$  implies  $I(C : B)_\sigma \leq \varepsilon$ .

Analogously to the exact case, unitary schemes satisfying approximate secrecy are equivalent to approximate one-designs (see the full version of this article [3]).

### 4.1 Approximate non-malleability

**Definition.** We now define a natural approximate-case analogue of NM, i.e., [Definition 3.4](#). Let us briefly recall the context. The malleability scenario is described by systems  $A, C, B$  and  $R$  (respectively, plaintext, ciphertext, side-information, and reference), an initial tripartite state  $\varrho_{ABR}$ , and an attack channel  $\Lambda_{CB \rightarrow C\tilde{B}}$ . Given this data, we have the effective channel  $\tilde{\Lambda}_{AB \rightarrow A\tilde{B}}$  defined in Equation (3.1) and the “unavoidable attack” probability  $p_{=}(\Lambda, \varrho)$  defined in Equation (3.2). The new definition now simply relaxes the requirement on the increase of the adversary’s mutual information.

**Definition 4.2 (Approximate non-malleability)** A QES  $(\tau_K, E, D)$  is  $\varepsilon$ -non-malleable ( $\varepsilon$ -NM) if for any state  $\varrho_{ABR}$  and any CPTP map  $\Lambda_{CB \rightarrow C\tilde{B}}$ , we have

$$I(AR : \tilde{B})_{\tilde{\Lambda}(\varrho)} \leq I(AR : B)_\varrho + h(p_{=}(\Lambda, \varrho)) + \varepsilon. \quad (4.1)$$

We record the approximate version of [Proposition 3.5](#), i.e., non-malleability implies secrecy. The proof is a straightforward adaptation of the exact case.

**Proposition 4.3** Let  $(\tau_K, E, D)$  be an  $\varepsilon$ -NM QES. Then  $(\tau_K, E, D)$  is  $2\varepsilon$ -ITS.

**Non-malleability with approximate designs.** Continuing as before, we now generalize the characterization theorems of non-malleability ([Theorem 3.7](#) and [Theorem 3.6](#)) to the approximate case.

**Theorem 4.4** *Let  $(\tau, E, D)$  be a QES with ciphertext dimension  $|C| = 2^m$  and  $r > 0$  a sufficiently large constant. Then the following holds:*

1. *If  $(\tau, E, D)$  is  $2^{-rm}$ -NM, then for any attack  $\Lambda_{CB \rightarrow C\bar{B}}$ , the effective map  $\tilde{\Lambda}_{AB \rightarrow A\bar{B}}$  is  $2^{-\Omega(m)}$ -close (in diamond norm) to*

$$\tilde{\Lambda}_{AB \rightarrow A\bar{B}}^{\text{exact}} = \text{id}_A \otimes \Lambda'_{B \rightarrow \bar{B}} + \frac{1}{|C|^2 - 1} (|C|^2 \langle D_K(\tau) \rangle - \text{id})_A \otimes \Lambda''_{B \rightarrow \bar{B}},$$

*with  $\Lambda', \Lambda''$  as in [Theorem 3.7](#).*

2. *Suppose that  $\log |R| = O(2^m)$ , where  $R$  is the reference register in [Definition 4.2](#). Then there exists a constant  $r$ , such that if every attack  $\Lambda_{CB \rightarrow C\bar{B}}$  results in an effective map that is  $2^{-rm}$ -close to  $\tilde{\Lambda}^{\text{exact}}$ , then the scheme is  $2^{-\Omega(m)}$ -NM.*

This theorem is proven with explicit constants in [Appendix B](#) as [Theorem B.3](#). The condition on  $R$  required for the second implication is necessary, as the relevant mutual information can at worst grow proportional to the logarithm of the dimension according to the Alicki-Fannes inequality. This is not a very strong requirement, as it should be relatively easy for the honest parties to put a bound on their total memory.

Next, we record the corollary which states that, for unitary schemes, approximate non-malleability is equivalent to encryption with an approximate 2-design. The proof proceeds as in the exact case, now starting from [Theorem 4.4](#).

**Theorem 4.5** *Let  $\Pi = (\tau_K, E, D)$  be a unitary QES for  $n$ -qubit messages and  $f : \mathbb{N} \rightarrow \mathbb{N}$  a function that grows at most exponential. Then there exists a constant  $r > 0$  such that*

1. *If  $\{E_k\}$  is a  $\Omega(2^{-rn})$ -approximate 2-design and  $\log |R| \leq f(n)$ , then  $\Pi$  is  $2^{-\Omega(n)}$ -NM.*
2. *If  $\Pi$  is  $\Omega(2^{-rn})$ -NM, then  $\{E_k\}_{k \in K}$  is a  $2^{-\Omega(n)}$ -approximate 2-design.*

**Relationship to approximate ABW.** Recall that, in [Section 3.2](#), we discussed the relationship between our notion of exact non-malleability and that of Ambainis et al. [\[6\]](#) (i.e., ABW-NM.) As we now briefly outline, our conclusions carry over to the approximate case without any significant changes.

As described in Equation (3'') of [\[6\]](#), one first relaxes the notion of ABW-NM appropriately by requiring that the containment [\(3.15\)](#) in [Definition 3.8](#) holds up to  $\varepsilon$  error in the diamond-norm distance. In the unitary case, both definitions are equivalent to approximate 2-designs (by the results of [\[6\]](#), and our [Theorem 4.5](#)). In the case of general schemes, the plaintext injection attack described in [Example 3.10](#) again shows that approximate ABW-NM is insufficient, and that approximate NM is strictly stronger.



## 4.2 Authentication

**Definitions.** Our definitions of authentication will be faithful to the original versions in [16, 18], with one slight modification. When decryption rejects, our encryption schemes (Definition 3.1) output  $\perp$  in the plaintext space, rather than setting an auxiliary qubit to a “reject” state. These definitions are equivalent in the sense that one can always set an extra qubit to “reject” conditioned on the plaintext being  $\perp$  (or vice-versa). Nonetheless, as we will see below, this mild change has some interesting consequences.

We begin with the definition of Dupuis, Nielsen and Salvail [16], which demands that the effective average channel of the attacker ignores the plaintext.

**Definition 4.6 (DNS Authentication [16])** *A QES  $(\tau_K, E, D)$  is called  $\varepsilon$ -DNS-authenticating if, for any CPTP-map  $\Lambda_{CB \rightarrow CB'}$ , there exists CP-maps  $\Lambda_{B \rightarrow \bar{B}}^{\text{acc}}$  and  $\Lambda_{B \rightarrow \bar{B}}^{\text{rej}}$  such that  $\Lambda^{\text{acc}} + \Lambda^{\text{rej}}$  is<sup>1</sup> TP, and for all  $\varrho_{AB}$  we have*

$$\left\| \text{Tr}_K D(\Lambda(E(\varrho_{AB} \otimes \tau_K))) - (\Lambda^{\text{acc}}(\varrho_{AB}) + |\perp\rangle\langle\perp| \otimes \Lambda^{\text{rej}}(\varrho_B)) \right\|_1 \leq \varepsilon. \quad (4.2)$$

An alternative definition was recently given by Garg, Yuen and Zhandry [18]. It asks that, *conditioned on acceptance*, with high probability the effective channel is close to a channel which ignores the plaintext.

**Definition 4.7 (GYZ Authentication [18])** *A QES  $(\tau_K, E, D)$  is called  $\varepsilon$ -GYZ-authenticating if, for any CPTP-map  $\Lambda_{CB \rightarrow CB'}$ , there exists a CP-map  $\Lambda_{B \rightarrow \bar{B}}^{\text{acc}}$  such that for all  $\varrho_{AB}$*

$$\left\| \Pi_{\text{acc}} D(\Lambda(E(\varrho_{AB} \otimes \tau_K))) \Pi_{\text{acc}} - \Lambda^{\text{acc}}(\varrho_{AB}) \otimes \tau_K \right\|_1 \leq \varepsilon. \quad (4.3)$$

Here  $\Pi_{\text{acc}}$  is the acceptance projector, i.e. projection onto  $\mathcal{H}_A$  in  $\mathcal{H}_A \oplus \mathbb{C}|\perp\rangle$ .

A peculiar aspect of the original definition in [18] is that it does not specify the outcome in case of rejection, and is thus stated in terms of trace non-increasing maps. Of course, all realistic quantum maps must be CPTP; this means that the designer of the encryption scheme must still declare what to do with the contents of the plaintext register after decryption. Our notion of decryption makes one such choice (i.e., output  $\perp$ ) which seems natural.

**GYZ authentication implies DNS authentication.** A priori, the relationship between Definition 2.2 in [16] and Definition 8 in [18] is not completely clear. On one hand, the latter is stronger in the sense that it requires success with high probability (rather than simply on average.) On the other hand, the former makes the additional demand that the ciphertext is untouched even if we reject. As we now show, GYZ-authentication in fact implies DNS-authentication.

<sup>1</sup> Note that there is a typographic error in [16] and [11] at this point of the definition. In those papers, the two effective maps are asked to sum to the identity (instead of just a TP map), which is impossible for many obvious choices of  $\Lambda$ .

**Theorem 4.8** *Let  $(\tau, E, D)$  be  $\varepsilon$ -totally authenticating for sufficiently small  $\varepsilon$ . Then it is  $O(\sqrt{\varepsilon})$ -DNS authenticating.*

*Proof.* Let  $\Lambda_{CB \rightarrow C\tilde{B}}$  be a CPTP map and  $\varepsilon \leq 62^{-2}$ . By Definition 4.7 there exists a CP map  $\Lambda'_{B \rightarrow \tilde{B}}$  such that for all states  $\varrho_{AB}$ ,

$$\| \Pi_a D(\Lambda(E(\varrho_{AB} \otimes \tau_K)) \Pi_a - \Lambda'(\varrho_{AB} \otimes \tau_K)) \|_1 \leq \varepsilon. \quad (4.4)$$

Assume for simplicity that  $D = M_\perp \circ D$ , where  $M_\perp$  measures the rejection symbol versus the rest. (otherwise we can define a new decryption map that way.) Define the CP maps

$$\begin{aligned} \Lambda_{AB \rightarrow \tilde{B}}^{(1)} &= \text{Tr}_A \Pi_a \tilde{\Lambda}(\cdot) & \Lambda_{AB \rightarrow \tilde{B}}^{(2)} &= \langle \perp |_A \tilde{\Lambda}(\cdot) | \perp \rangle_A \\ \Lambda''_{B \rightarrow \tilde{B}} &= \text{Tr}_C \Lambda(E_K(\tau_A) \otimes (\cdot)). \end{aligned}$$

By Theorem 15 in [18] we have

$$\| E_K(\varrho_{ABR}) - E_K(\tau_A) \otimes \varrho_{BR} \|_1 \leq 14\sqrt{\varepsilon}, \quad (4.5)$$

which implies that

$$\| \text{Tr}_A \otimes \Lambda'' - \text{Tr}_C \circ \Lambda \circ E_K \|_\diamond \leq \hat{\varepsilon} := 14\sqrt{\varepsilon}. \quad (4.6)$$

Note that

$$\begin{aligned} \text{Tr}_C \circ \Lambda \circ E_K &= \text{Tr}_{CK} \circ \Lambda \circ E((\cdot) \otimes \tau_K) \\ &= \text{Tr}_{AK} \circ D \circ \Lambda \circ E((\cdot) \otimes \tau_K) = \text{Tr}_A \circ \tilde{\Lambda}. \end{aligned} \quad (4.7)$$

On the other hand, we also have that, by Equation (4.4),

$$\| \text{Tr}_A \circ \tilde{\Lambda} - \text{Tr}_A \otimes \Lambda' - \Lambda^{(2)} \| \leq \| \text{Tr}_A (\Pi_a \tilde{\Lambda}(\cdot)) - \Lambda' \|_\diamond \leq \varepsilon \quad (4.8)$$

Combining Equations (4.6), (4.7) and (4.8), we get

$$\| \Lambda^{(2)} - \text{Tr}_A \otimes (\Lambda'' - \Lambda') \|_\diamond \leq \varepsilon + \hat{\varepsilon}. \quad (4.9)$$

Now observe that

$$[\text{Tr}_A \otimes (\Lambda' - \Lambda'')]_{B \rightarrow \tilde{B}} \circ \Xi_{A \rightarrow A} = \text{Tr}_A \otimes (\Lambda' - \Lambda'')_{B \rightarrow \tilde{B}} \quad (4.10)$$

For all CPTP maps  $\Xi_{A \rightarrow A}$ . We define  $\Lambda'''_{B \rightarrow \tilde{B}} = \Lambda^{(2)}(\tau_A \otimes (\cdot))$  and calculate

$$\begin{aligned} \| \Lambda^{(2)} - \text{Tr}_A \otimes \Lambda''' \|_\diamond &\leq \| \Lambda^{(2)} - \text{Tr}_A \otimes (\Lambda'' - \Lambda') \|_\diamond \\ &\quad + \| \text{Tr}_A \otimes (\Lambda'' - \Lambda') - \text{Tr}_A \otimes \Lambda''' \|_\diamond, \end{aligned}$$

by the triangle inequality for the diamond norm. Continuing with the calculation,

$$\begin{aligned} \| \Lambda^{(2)} - \text{Tr}_A \otimes \Lambda''' \|_\diamond &\leq \varepsilon + \hat{\varepsilon} + \| \text{Tr}_A \otimes (\Lambda'' - \Lambda') - \text{Tr}_A \otimes \Lambda''' \|_\diamond \\ &= \varepsilon + \hat{\varepsilon} + \| [\text{Tr}_A \otimes (\Lambda'' - \Lambda') - \Lambda^{(2)}] \circ \langle \tau_A \rangle_{A \rightarrow A} \|_\diamond \\ &\leq 2(\varepsilon + \hat{\varepsilon}) = 28\sqrt{\varepsilon} + 2\varepsilon. \end{aligned} \quad (4.11)$$

The first inequality above is Equation (4.9). The first equality is just a rewriting of the definition of  $\Lambda'''$ , and the second equality is Equation (4.10). Finally, the last inequality is due to Equation (4.9) and the fact that the diamond norm is submultiplicative.

We have almost proven security according to Definition 4.6, as we have shown  $\tilde{\Lambda}$  to be close in diamond norm to  $\text{id}_A \otimes \Lambda' + \langle |\perp\rangle\langle\perp| \rangle \otimes \Lambda'''$ . However,  $\Lambda' + \Lambda'''$  is only approximately TP; more precisely, we have that for all  $\varrho_{ABR}$ ,

$$|\text{Tr}(\Lambda' + \Lambda''')(\varrho_{ABR}) - 1| \leq 28\sqrt{\varepsilon} + 3\varepsilon \quad (4.12)$$

by the triangle inequality. We therefore have to modify  $\Lambda' + \Lambda'''$  so that it becomes TP, while keeping the structure required for DNS authentication. Let  $M_B = (\Lambda' + \Lambda''')^\dagger(\mathbb{1}_{\tilde{B}})$ . (4.12). Defining the CP-map  $\mathcal{M}(X) = M^{-1/2}XM^{-1/2}$  and noting it is well-behaved for small  $\varepsilon$ , it follows from a straightforward calculation (see the full version [3] of this article for details) that

$$\left\| \tilde{\Lambda}_{AB \rightarrow A\tilde{B}} - \text{id}_A \otimes \Lambda_{B \rightarrow \tilde{B}}^{\text{acc}} - \perp \otimes \Lambda_{B \rightarrow \tilde{B}}^{\text{rej}} \right\|_{\diamond} \leq O(\sqrt{\varepsilon}). \quad (4.13)$$

with  $\lambda^{\text{acc}} = \Lambda' \circ \mathcal{M}$  and  $\Lambda^{\text{rej}} = \Lambda'' \circ \mathcal{M}$ .

□

**Achieving GYZ authentication with two-designs.** In [18], the authors provide a scheme for their notion of authentication based on unitary eight-designs. We now show that, in fact, an approximate 2-design suffices. This implies that the well-known Clifford scheme (see e.g [15, 11]) satisfies the strong security of Definition 4.7. We remark that our proof is inspired by the reasoning based on Schur's lemma used in results on decoupling [8, 17, 24, 9].

**Theorem 4.9** *Let  $D = \{U_k\}_k$  be a  $\delta$ -approximate unitary 2-design on  $\mathcal{H}_C$ . Let  $\mathcal{H}_C = \mathcal{H}_A \otimes \mathcal{H}_T$  and define*

$$\begin{aligned} E_k(X_A) &= U_k(X_A \otimes |0\rangle\langle 0|_T)(U_k)^\dagger \\ D_k(Y_C) &= \langle 0|_T(U_k)^\dagger Y U_k |0\rangle_T + \text{Tr}((\mathbb{1}_T - |0\rangle\langle 0|_T)(U_k)^\dagger Y U_k) |\perp\rangle\langle\perp|. \end{aligned}$$

*Then the QES  $(\tau_K, E, D)$  is  $4(1/|T| + 3\delta)^{1/3}$ -GYZ-authenticating.*

**Remark 4.10** *The following proof uses the same simulator as the proof for the 8-design scheme in [18], called "oblivious adversary" there. The construction exhibited there is efficient given that the real adversary is efficient.*

*Proof.* To improve readability, we will occasionally switch between adding subscripts to operators (indicating which spaces they act on) and omitting these subscripts. We begin by remarking that it is sufficient to prove the GYZ condition (specifically, Equation 4.3) for pure input states and isometric adversary channels. Indeed, for a general state  $\varrho_{AB}$  and a general map  $\Lambda_{CB \rightarrow C\tilde{B}}$ , we may

let  $\varrho_{ABR}$  and  $V_{CB \rightarrow C\bar{B}E}$  be the purification and Stinespring dilation, respectively. We then simply observe that the trace distance decreases under partial trace (see e.g. [25]). Let  $\varrho_{AB}$  be a pure input state and

$$\Lambda_{CB \rightarrow C\bar{B}}(X_{CB}) = V_{CB \rightarrow C\bar{B}} X_{CB} V_{CB \rightarrow C\bar{B}}^\dagger$$

an isometry. We define the corresponding ‘‘ideal’’ channel  $\Gamma_V$ , and the corresponding ‘‘real, accept’’ channel  $\Phi_k$ , as follows:

$$\begin{aligned} (\Gamma_V)_{B \rightarrow \bar{B}} &= \frac{1}{|C|} \text{Tr}_C V \text{ and} \\ (\Phi_k)_{AB \rightarrow A\bar{B}} &= \langle 0|_T (U_k)_C^\dagger V_{CB \rightarrow C\bar{B}} U_k |0\rangle_T. \end{aligned} \quad (4.14)$$

Note that for any matrix  $M$  with  $\|M\|_\infty \leq 1$ , the map  $\Lambda_M(X) = M^\dagger X M$  is completely positive and trace non-increasing. We have

$$\|\Gamma_V\|_\infty \leq \frac{1}{|C|} \sum_i \|\langle i|V|i\rangle\|_\infty \leq 1. \quad (4.15)$$

We start by bounding the expectation of  $\|((\Gamma_V)_{B \rightarrow \bar{B}} - (\Phi_k)_{AB \rightarrow A\bar{B}})|\varrho\rangle_{AB}\|_2^2$ , as follows. To simplify notation, we set  $\sigma_{ABT} := |\varrho\rangle\langle\varrho|_{AB} \otimes |0\rangle\langle 0|_T$  to be the tagged state corresponding to plaintext (and side information)  $\varrho_{AB}$ .

$$\begin{aligned} \frac{1}{|K|} \sum_k \|(\Gamma_V - \Phi_k)|\varrho\rangle\|_2^2 &= \frac{1}{|K|} \sum_k \langle\varrho|(\Gamma_V - \Phi_k)^\dagger(\Gamma_V - \Phi_k)|\varrho\rangle \\ &= \frac{1}{|K|} \sum_k \text{Tr} [\sigma_{ABT}(U_k)^\dagger V^\dagger U_k |0\rangle\langle 0| (U_k)^\dagger V U_k] \\ &\quad - 2 \frac{1}{|K|} \sum_k \text{Tr} [\sigma_{ABT}(U_k)^\dagger V^\dagger U_k \Gamma_V] + \langle\varrho|(\Gamma_V)^\dagger \Gamma_V|\varrho\rangle. \end{aligned} \quad (4.16)$$

First we bound the second term, using the fact that  $\Gamma_V$  only acts on  $B$ .

$$\begin{aligned} \frac{1}{|K|} \sum_k \text{Tr} [\sigma_{ABT}(U_k)^\dagger V^\dagger U_k \Gamma_V] &= \frac{1}{|K|} \sum_k \text{Tr} [U_k \sigma_{ABT}(U_k)^\dagger V^\dagger \Gamma_V] \\ &= \int \text{Tr} [(U \sigma_{ABT} U^\dagger + \Delta) V^\dagger \Gamma_V] \geq \int \text{Tr} [U \sigma_{ABT} U^\dagger V^\dagger \Gamma_V] - \delta \\ &= \int \text{Tr} [\sigma_{ABT} U^\dagger V^\dagger U \Gamma_V] - \delta = \langle\varrho|(\Gamma_V)^\dagger \Gamma_V|\varrho\rangle - \delta. \end{aligned} \quad (4.17)$$

In the above, the operator  $\Delta$  is the ‘‘error’’ operator in the  $\delta$ -approximate 2-design. The second equality above follows from  $\|\Delta\|_1 \leq \delta$  and the fact that a 2-design is also a 1-design; the inequality follows by Hölder’s inequality, and the last step follows from Schur’s lemma.

The first term of the RHS of Equation (4.16) can be simplified as follows. We will begin by applying the swap trick (Lemma 2.1)  $\text{Tr}[XY] = \text{Tr}[FX \otimes Y]$  in the

second line below. The swap trick is applied to register  $CC'$ , with the operators  $X$  and  $Y$  defined as indicated below.

$$\begin{aligned}
& \frac{1}{|K|} \sum_k \text{Tr} \left[ \underbrace{\sigma_{ABT}(U_k)_C^\dagger V_{C\bar{B} \rightarrow CB}^\dagger (U_k)_C |0\rangle\langle 0|_T}_{X} \underbrace{(U_k)_C^\dagger V_{CB \rightarrow C\bar{B}}(U_k)_C}_{Y} \right] \\
&= \frac{1}{|K|} \sum_k \text{Tr} \left[ (\sigma_{ABT} \otimes |0\rangle\langle 0|_{T'}) (U_k^{\otimes 2})_{CC'} V_{C\bar{B} \rightarrow CB}^\dagger V_{C'B \rightarrow C'\bar{B}} (U_k^{\otimes 2})_{CC'}^\dagger F_{CC'} \right] \\
&= \frac{1}{|K|} \sum_k \text{Tr} \left[ (U_k^{\otimes 2})_{CC'}^\dagger (\sigma_{ABT} \otimes |0\rangle\langle 0|_{T'}) (U_k^{\otimes 2})_{CC'} V_{C\bar{B} \rightarrow CB}^\dagger V_{C'B \rightarrow C'\bar{B}} F_{CC'} \right] \\
&\leq \int \text{Tr} \left[ (U^{\otimes 2})_{CC'}^\dagger (\sigma_{ABT} \otimes |0\rangle\langle 0|_{T'}) U_{CC'}^{\otimes 2} V_{C\bar{B} \rightarrow CB}^\dagger V_{C'B \rightarrow C'\bar{B}} F_{CC'} \right] + \delta \\
&= \int \text{Tr} \left[ (\sigma_{ABT} \otimes |0\rangle\langle 0|_{T'}) U_{CC'}^{\otimes 2} V_{C\bar{B} \rightarrow CB}^\dagger V_{C'B \rightarrow C'\bar{B}} (U^{\otimes 2})_{CC'}^\dagger F_{CC'} \right] + \delta.
\end{aligned} \tag{4.18}$$

The inequality above follows the same way as in [Equation 4.17](#). Let  $d = |C|$ . An easy representation-theoretic calculation (see the Full version [3] for details) shows that

$$\int U^{\otimes 2} V_{C\bar{B} \rightarrow CB}^\dagger V_{C'B \rightarrow C'\bar{B}} (U^{\otimes 2})^\dagger dU = \mathbb{1}_{CC'} \otimes R_B^1 + F_{CC'} \otimes R_B^F, \tag{4.19}$$

where we have set

$$\begin{aligned}
R_B^1 &= \frac{1}{d(d^2-1)} \left( d^3 \Gamma_V^\dagger \Gamma_V - d\mathbb{1} \right) = \frac{1}{(d^2-1)} \left( d^2 \Gamma_V^\dagger \Gamma_V - \mathbb{1} \right) \\
R_B^F &= \frac{1}{d(d^2-1)} \left( d^2 \mathbb{1} - d^2 \Gamma_V^\dagger \Gamma_V \right) = \frac{d}{(d^2-1)} \left( \mathbb{1} - \Gamma_V^\dagger \Gamma_V \right).
\end{aligned}$$

plugging [\(4.19\)](#) into [\(4.18\)](#) and using [Lemma 2.1](#) again, we get

$$\begin{aligned}
& \int \text{Tr} \left[ (\sigma_{ABT} \otimes |0\rangle\langle 0|_{T'}) U_{CC'}^{\otimes 2} V_{C\bar{B} \rightarrow CB}^\dagger V_{C'B \rightarrow C'\bar{B}} (U^{\otimes 2})_{CC'}^\dagger F_{CC'} \right] \\
&= \text{Tr} \left[ (\sigma_{ABT} \otimes |0\rangle\langle 0|_{T'}) (\mathbb{1}_{CC'} \otimes R_{B^2 \rightarrow \bar{B}^2}^1 + F_{CC'} \otimes R_{B^2 \rightarrow \bar{B}^2}^F) F_{CC'} \right] \\
&= \text{Tr} \left[ |\varrho\rangle\langle \varrho|_B (R_B^1 + |A|R_B^F) \right] \\
&= \text{Tr} \left[ |\varrho\rangle\langle \varrho|_B \left( \frac{d(d-|A|)}{d^2-1} (\Gamma_V^\dagger \Gamma_V)_B + \frac{d|A|-1}{d^2-1} \mathbb{1}_B \right) \right].
\end{aligned} \tag{4.20}$$

Now recall that  $d = |A||T|$ . Using the fact that  $(a-1)/(b-1) \leq a/b$  for  $b \geq a$ , we can give a bound as follows.

$$\begin{aligned}
& \text{Tr} \left[ |\varrho\rangle\langle \varrho| \left( \frac{d(d-|A|)}{d^2-1} (\Gamma_V^\dagger \Gamma_V) + \frac{d|A|-1}{d^2-1} \mathbb{1} \right) \right] \\
&= \frac{d|A|(|T|-1)}{d^2-1} \langle \varrho | (\Gamma_V^\dagger \Gamma_V) | \varrho \rangle + \frac{d|A|-1}{d^2-1} \\
&\leq \langle \varrho | (\Gamma_V^\dagger \Gamma_V) | \varrho \rangle + \frac{1}{|T|}.
\end{aligned} \tag{4.21}$$

Putting everything together, we arrive at

$$\frac{1}{|K|} \sum_k \|(\Gamma_V - \Phi_k)|\varrho\rangle\|_2^2 \leq \frac{1}{|T|} + 3\delta. \quad (4.22)$$

By Markov's inequality this implies

$$\mathbb{P} \left[ \|(\Gamma_V - \Phi_k)|\varrho\rangle\|_2^2 > \alpha \left( \frac{1}{|T|} + 3\delta \right) \right] \leq \frac{1}{\alpha} \quad (4.23)$$

which is equivalent to

$$\mathbb{P} \left[ \|(\Gamma_V - \Phi_k)|\varrho\rangle\|_2 > \alpha^{1/2} \left( \frac{1}{|T|} + 3\delta \right)^{1/2} \right] \leq \frac{1}{\alpha}, \quad (4.24)$$

where the probability is taken over the uniform distribution on  $\mathcal{D}$ . Choosing  $\alpha = (1/|T| + 3\delta)^{-1/3}$  this yields

$$\mathbb{P} \left[ \|(\Gamma_V - \Phi_k)|\varrho\rangle\|_2 > \left( \frac{1}{|T|} + 3\delta \right)^{1/3} \right] \leq \left( \frac{1}{|T|} + 3\delta \right)^{1/3}. \quad (4.25)$$

Let  $S \subset \mathcal{D}$  be such that  $|S|/|\mathcal{D}| \geq 1 - (1/|T| + 3\delta)^{1/3}$  and  $\|(\Gamma_V - \Phi_k)|\varrho\rangle\|_2 \leq (1/|T| + 3\delta)^{1/3}$  for all  $U_k \in S$ . Using the easy-to-verify inequality  $\|\psi\rangle\langle\psi| - |\phi\rangle\langle\phi|\|_1 \leq 2\|\psi\rangle - |\phi\rangle\|_2$ , we can bound

$$\begin{aligned} & \frac{1}{|K|} \sum_{U_k \in \mathcal{D}} \left\| \Phi_k |\varrho\rangle\langle\varrho| (\Phi_k)^\dagger - \Gamma_V |\varrho\rangle\langle\varrho| \Gamma_V^\dagger \right\|_1 \\ & \leq \frac{1}{|S|} \sum_{U_k \in S} \left\| \Phi_k |\varrho\rangle\langle\varrho| (\Phi_k)^\dagger - \Gamma_V |\varrho\rangle\langle\varrho| \Gamma_V^\dagger \right\|_1 + 2 \left( \frac{1}{|T|} + 3\delta \right)^{1/3} \\ & \leq \frac{2}{|S|} \sum_{U_k \in S} \|(\Gamma_V - \Phi_k)|\varrho\rangle\|_2 + 2|T|^{-1/3} \\ & \leq 4 \left( \frac{1}{|T|} + 3\delta \right)^{1/3}. \end{aligned} \quad (4.26)$$

This completes the proof for pure states and isometric adversary channels. As noted above, the general case follows.  $\square$

As an example, one may set  $|T| = 2^s$  (i.e.  $s$  tag qubits) and take an approximate unitary 2-design of accuracy  $2^{-s}$ . The resulting scheme would then be  $\Omega(2^{-s/3})$ -GYZ-authenticating.

A straightforward corollary of the above result is that, in the case of unitary schemes, adding tags to non-malleable schemes results in GYZ authentication. We leave open the question of whether this is the case for general (not necessarily unitary) schemes.

**Corollary 4.11** *Let  $(\tau, E, D)$  be a  $2^{-rn}$ -non-malleable unitary QES with plaintext space  $A$ . Define a new scheme  $(\tau, E', D')$  with plaintext space  $A'$  where  $A = TA'$  and*

$$\begin{aligned} E'(X) &= E(X \otimes |0\rangle\langle 0|_T) \\ D'(Y) &= \langle 0|_T D(Y) |0\rangle_T + \text{Tr}[(\mathbf{1}_T - |0\rangle\langle 0|_T) D(Y)] |\perp\rangle\langle \perp|. \end{aligned}$$

*Then there is a constant  $r > 0$  such that  $(\tau, E', D')$  is  $2^{-\Omega(n)}$ -GYZ-authenticating if  $|T| = 2^{\Omega(n)}$ .*

The proof is a direct application of [Theorem 4.5](#) (approximate non-malleability is equivalent to approximate 2-design) and [Theorem 4.9](#) (approximate 2-designs suffice for GYZ authentication.) We emphasize that, by [Remark 2.3](#), exponential accuracy requirements can be met with polynomial-size circuits.

**DNS authentication from non-malleability.** We end with a theorem concerning the case of general (i.e., not necessarily unitary) schemes. We show that adding tags to a non-malleable scheme results in a DNS-authenticating scheme. In this proof we will denote the output system of the decryption map by  $\bar{A}$  to emphasize that it is  $A$  enlarged by the reject symbol.

**Theorem 4.12** *Let  $r$  be a sufficiently large constant, and let  $(\tau, E, D)$  be an  $2^{-rn}$ -NM QES with  $n$  qubit plaintext space  $A$ , and choose an integer  $d$  dividing  $|A|$ . Then there exists a decomposition  $A = TA'$  and a state  $|\psi\rangle_T$  such that  $|T| = d$  and the scheme  $(\tau, E', D')$  defined by*

$$\begin{aligned} E^t(X) &= E(X \otimes |\psi\rangle\langle \psi|_T) \\ D^t(Y) &= \langle \psi|_T D(Y) |\psi\rangle_T + \text{Tr}[(\mathbf{1}_T - |\psi\rangle\langle \psi|_T) D(Y)] |\perp\rangle\langle \perp|. \end{aligned}$$

*is  $(4/|T|) + 2^{-\Omega(n)}$ -DNS-authenticating.*

*Proof.* We prove the statement for  $\varepsilon = 0$  for simplicity, the general case follows easily by employing [Theorem 4.4](#) instead of [Theorem 3.7](#).

By [Theorem 3.7](#), for any attack map  $\Lambda_{CB \rightarrow C\bar{B}}$ , the effective map is equal to

$$\tilde{\Lambda}_{AB \rightarrow \bar{A}\bar{B}} = \text{id}_A \otimes \Lambda'_{B \rightarrow \bar{B}} + \frac{1}{|C|^2 - 1} (|C|^2 \langle D_K(\tau_C) \rangle - \text{id})_{\bar{A}} \otimes \Lambda''_{B \rightarrow \bar{B}} \quad (4.27)$$

for CP maps  $\Lambda'$  and  $\Lambda''$  whose sum is TP. The effective map under the tagged scheme is therefore

$$\begin{aligned} \tilde{\Lambda}_{A'B \rightarrow \bar{A}'\bar{B}}^t &= \langle \psi|_T \tilde{\Lambda}_{AB \rightarrow \bar{A}\bar{B}}((\cdot) \otimes \psi_T) |\psi\rangle_T \\ &\quad + \text{Tr}[(\mathbf{1}_T - \psi_T) \tilde{\Lambda}_{AB \rightarrow \bar{A}\bar{B}}((\cdot) \otimes \psi_T)] |\perp\rangle\langle \perp| \\ &= (\text{id}_{A'})_{A' \rightarrow \bar{A}'} \otimes \Lambda'_{B \rightarrow \bar{B}} \\ &\quad + (|C|^2 \langle (\langle \psi|_T D_K(\tau_C) |\psi\rangle_T)_{A'} \oplus \beta |\perp\rangle\langle \perp| \rangle - \text{id}_{A'})_{A \rightarrow \bar{A}'} \otimes \frac{\Lambda''_{B \rightarrow \bar{B}}}{|C|^2 - 1} \end{aligned}$$

with  $\beta = \text{Tr}[(\mathbb{1} - \psi)_T D_K(\tau_C)]$ . We would like to say that, unless the output is the reject symbol, the effective map on  $A$  is the identity. We do not know, however, what  $D_K(\tau_C)$  looks like. Therefore we apply a standard reasoning that if a quantity is small *in expectation*, then there exists at least one small instance. We calculate the expectation of  $\text{Tr}\langle\psi|_T D_K(\tau_C)|\psi\rangle_T$  when the decomposition  $A = TA'$  is drawn at random according to the Haar measure,

$$\begin{aligned} \int \text{Tr}\langle\psi|U_A^\dagger D_K(\tau_C)U_A|\psi\rangle_T dU_A &= \text{Tr}\left[\left(\int U_A|\psi\rangle_T \otimes \mathbb{1}_{A'}\psi U_A^\dagger dU_A\right) D_K(\tau_C)\right] \\ &= \frac{\text{Tr}\mathbb{1}_A}{\text{Tr}\Pi_{\text{acc}}} \text{Tr}\Pi_{\text{acc}} D_K(\tau_C) \\ &\leq 1/|T|. \end{aligned} \tag{4.28}$$

Hence there exists at least one decomposition  $A = TA'$  and a state  $|\psi\rangle_T$  such that  $\hat{\gamma} := \text{Tr}\langle\psi|_T D_K(\tau_C)|\psi\rangle_T \leq 1/|T|$ . Define  $\gamma = \max(\hat{\gamma}, |C|^{-2})$ . For the resulting primed scheme, let

$$A_{\text{rej}} := \frac{(1 - \gamma)|C|^2}{|C|^2 - 1} A'' \quad \text{and} \quad A_{\text{acc}} = A' + \frac{\gamma|C|^2 - 1}{|C|^2 - 1} A''.$$

We calculate the diamond norm difference between the real effective map and the ideal effective map,

$$\begin{aligned} &\|\tilde{A}^t - \text{id} \otimes A_{\text{acc}} - \langle\perp|\rangle\langle\perp| \otimes A_{\text{rej}}\|_{\diamond} \\ &\leq \|\text{id} \otimes A' + \frac{1}{|C|^2 - 1} (|C|^2 \langle\langle\psi|D_K(\tau)|\psi\rangle\rangle - \text{id}) \otimes A'' - \text{id} \otimes A_{\text{acc}}\|_{\diamond} \\ &\quad + \|\langle\perp|\rangle\langle\perp| \otimes (1 - \hat{\gamma})|C|^2 A'' / (|C|^2 - 1) - \langle\perp|\rangle\langle\perp| \otimes A_{\text{rej}}\|_{\diamond} \\ &\leq (1 + |C|^{-2})(|T|^{-1} + 2|C|^{-2}) \\ &= |T|^{-1}(1 + (|A'| |T|)^{-2})(1 + 2|A'|^{-2}) \\ &\leq 4|T|^{-1} \end{aligned} \tag{4.29}$$

as desired.  $\square$

## Acknowledgments

The authors would like to thank Anne Broadbent, Alexander Müller-Hermes, Frédéric Dupuis and Christopher Portmann for helpful discussions. G.A. and C.M. acknowledge financial support from the European Research Council (ERC Grant Agreement 337603), the Danish Council for Independent Research (Sapere Aude) and VILLUM FONDEN via the QMATH Centre of Excellence (Grant 10059).



## Bibliography

- [1] Scott Aaronson and Daniel Gottesman. Improved simulation of stabilizer circuits. *Phys. Rev. A*, 70:052328, Nov 2004. doi: 10.1103/PhysRevA.70.052328.
- [2] Dorit Aharonov, Michael Ben-Or, and Elad Eban. Interactive proofs for quantum computations. In *Innovations in Computer Science - ICS 2010, Tsinghua University, Beijing, China, January 5-7, 2010. Proceedings*, pages 453–469, 2010.
- [3] Gorjan Alagic and Christian Majenz. Quantum non-malleability and authentication. *CoRR*, abs/1610.04214, 2016. URL <http://arxiv.org/abs/1610.04214>.
- [4] Robert Alicki and Mark Fannes. Continuity of quantum conditional information. *Journal of Physics A: Mathematical and General*, 37(5):L55, 2004.
- [5] Andris Ambainis, Michele Mosca, Alain Tapp, and Ronald De Wolf. Private quantum channels. In *focs*, pages 547–553, 2000.
- [6] Andris Ambainis, Jan Bouda, and Andreas Winter. Nonmalleable encryption of quantum information. *Journal of Mathematical Physics*, 50(4):042106, 2009.
- [7] Howard Barnum, Claude Crépeau, Daniel Gottesman, Adam Smith, and Alain Tapp. Authentication of quantum messages. In *Foundations of Computer Science, 2002. Proceedings. The 43rd Annual IEEE Symposium on*, pages 449–458. IEEE, 2002.
- [8] Mario Berta, Matthias Christandl, and Renato Renner. The quantum reverse shannon theorem based on one-shot information theory. *Communications in Mathematical Physics*, 306(3):579–615, 2011.
- [9] Mario Berta, Fernando GSL Brandao, Christian Majenz, and Mark M Wilde. Deconstruction and conditional erasure of quantum correlations. *arXiv preprint arXiv:1609.06994*, 2016.
- [10] Fernando GSL Brandao, Aram W Harrow, and Michal Horodecki. Local random quantum circuits are approximate polynomial-designs. *arXiv preprint arXiv:1208.0692*, 2012.
- [11] Anne Broadbent and Evelyn Wainwright. Efficient simulation for quantum message authentication. *arXiv preprint arXiv:1607.03075*, 2016.
- [12] Man-Duen Choi. Completely positive linear maps on complex matrices. *Linear algebra and its applications*, 10(3):285–290, 1975.
- [13] Richard Cleve, Debbie Leung, Li Liu, and Chunhao Wang. Near-linear constructions of exact unitary 2-designs. *Quantum Information and Computation*, 16(9&10):0721–0756, 2016.
- [14] Christoph Dankert, Richard Cleve, Joseph Emerson, and Etera Livine. Exact and approximate unitary 2-designs and their application to fidelity estimation. *Physical Review A*, 80(1):012304, 2009.

- [15] Frédéric Dupuis, Jesper Buus Nielsen, and Louis Salvail. Secure two-party quantum evaluation of unitaries against specious adversaries. In *Annual Cryptology Conference*, pages 685–706. Springer, 2010.
- [16] Frédéric Dupuis, Jesper Buus Nielsen, and Louis Salvail. Actively secure two-party evaluation of any quantum operation. In *Advances in Cryptology–CRYPTO 2012*, pages 794–811. Springer, 2012.
- [17] Frédéric Dupuis, Mario Berta, Jürg Wullschleger, and Renato Renner. One-shot decoupling. *Communications in Mathematical Physics*, 328(1):251–284, 2014.
- [18] Sumegha Garg, Henry Yuen, and Mark Zhandry. New security notions and feasibility results for authentication of quantum data. *arXiv preprint arXiv:1607.07759*, 2016.
- [19] Andrzej Jamiołkowski. Linear transformations which preserve trace and positive semidefiniteness of operators. *Reports on Mathematical Physics*, 3(4):275–278, 1972.
- [20] Akinori Kawachi, Christopher Portmann, and Keisuke Tanaka. Characterization of the relations between information-theoretic non-malleability, secrecy, and authenticity. In *International Conference on Information Theoretic Security*, pages 6–24. Springer, 2011.
- [21] Elliott H Lieb and Mary Beth Ruskai. A fundamental property of quantum-mechanical entropy. *Physical Review Letters*, 30(10):434, 1973.
- [22] Elliott H. Lieb and Mary Beth Ruskai. Proof of the strong subadditivity of quantum-mechanical entropy. *Journal of Mathematical Physics*, 14(12):1938–1941, 1973.
- [23] Richard A Low. Pseudo-randomness and learning in quantum computation. *arXiv preprint arXiv:1006.5227*, 2010.
- [24] Christian Majenz, Mario Berta, Frédéric Dupuis, Renato Renner, and Matthias Christandl. Catalytic decoupling of quantum information. *arXiv preprint arXiv:1605.00514*, 2016.
- [25] Michael A Nielsen and Isaac L Chuang. *Quantum computation and quantum information*. Cambridge university press, 2010.
- [26] C. Portmann. Quantum authentication with key recycling. *ArXiv e-prints*, October 2016.
- [27] W Forrest Stinespring. Positive functions on  $c^*$ -algebras. *Proceedings of the American Mathematical Society*, 6(2):211–216, 1955.

## A Technical lemmas

In the following we state some technical Lemmas that we need in this article. The proofs can be found in the full version [3].

**Lemma A.1** *Let  $X_{A \rightarrow B} \in L(\mathcal{H}_A, \mathcal{H}_B)$  be a linear operator from  $A$  to  $B$ . Then*

$$X_{A \rightarrow B} |\phi^+\rangle_{AA'} = \sqrt{\frac{|B|}{|A|}} X_{B' \rightarrow A'}^T |\phi^+\rangle_{BB'}. \quad (\text{A.1})$$

The next group of lemmas is concerned with entropic quantities.

**Lemma A.2** Let  $\Lambda_{A \rightarrow A'}^{(i)}$  be CPTP maps and  $\Lambda_{B \rightarrow B'}^{(i)}$ ,  $i = 1, \dots, k$  CP maps for  $i = 1, \dots, k$  such that  $\sum_i \Lambda_{B \rightarrow B'}^{(i)}$  is trace preserving. Let  $\Lambda_{AB \rightarrow A'B'}^{(i)} = \Lambda_{A \rightarrow A'}^{(i)} \otimes \Lambda_{B \rightarrow B'}^{(i)}$  and define the CPTP maps

$$\begin{aligned} \Lambda_{AB \rightarrow A'B'C} &= \sum_{i=1}^k \Lambda_{AB \rightarrow A'B'}^{(i)} \otimes |i\rangle\langle i|_C \text{ and} \\ \Lambda'_{B \rightarrow B'C} &= \sum_{i=1}^k \Lambda_{B \rightarrow B'}^{(i)} \otimes |i\rangle\langle i|_C. \end{aligned} \quad (\text{A.2})$$

Then

$$I(A' : B')_{\Lambda(\varrho)} \leq I(A : B)_{\varrho} + H(C|A)_{\Lambda'(\varrho)} \leq I(A : B)_{\varrho} + H(C)_{\Lambda(\varrho)} \quad (\text{A.3})$$

for any quantum state  $\varrho_{AB}$ .

The final lemma characterizes CPTP maps that are invertible on their image such that the inverse is CPTP as well.

**Lemma A.3** Let  $(\tau_K, E, D)$  be a QES. Then the encryption maps have the structure

$$(E_k)_{A \rightarrow C} = (V_k)_{A\hat{C} \rightarrow C} \left( (\cdot) \otimes \sigma_{\hat{C}}^{(k)} \right) (V_k)_{A\hat{C} \rightarrow C}^\dagger, \quad (\text{A.4})$$

and the decryption maps hence must have the form

$$\begin{aligned} (D_k)_{C \rightarrow A} &= \text{Tr}_{\hat{C}} \left[ \Pi_{\text{supp}\sigma^k} (V_k)_{A\hat{C} \rightarrow C}^\dagger (\cdot) (V_k)_{A\hat{C} \rightarrow C} \right] \\ &+ \left( \hat{D}_k \right)_{C \rightarrow A} \left[ (\mathbb{1}_C - \Pi_k^{\text{valid}}) (\cdot) (\mathbb{1}_C - \Pi_k^{\text{valid}}) \right] \end{aligned} \quad (\text{A.5})$$

for some quantum states  $\sigma_{\hat{C}}^{(k)}$ , isometries  $(V_k)_{C \rightarrow A\hat{C}}$ , and some CPTP map  $\hat{D}_k$ . Here,  $\Pi_k^{\text{valid}} = (V_k)_{A\hat{C} \rightarrow C} \Pi_{\text{supp}\sigma^k} (V_k)_{A\hat{C} \rightarrow C}^\dagger$  is the projector onto the space of valid ciphertexts.

## B Proof of characterization theorem

This section is dedicated to proving the characterization theorem for non-malleable QES, i.e., [Theorem 4.4](#). We begin with two preparatory lemmas.

**Lemma B.1** For any QES  $(\tau, E, D)$  the map  $\mathcal{E} := |K|^{-1} \sum_k D_k \otimes E_k^T$  satisfies

$$\mathcal{E}(|\phi^+\rangle\langle\phi^+|_{CC'}(X \otimes \text{id}_{C'})) = \frac{|A|}{|C|} |\phi^+\rangle\langle\phi^+|_{AA'}(E_K^\dagger(X) \otimes \text{id}_{A'})$$

This lemma is a consequence of the correctness condition and is proven in the full version [\[3\]](#) of this article.

**Lemma B.2** Suppose  $(\tau_K, E, D)$  satisfies [Definition 4.2](#) for trivial  $B$ . Then  $\mathcal{E} := |K|^{-1} \sum_k D_k \otimes E_k^T$  satisfies

$$\begin{aligned} & \left\| \mathcal{E}(X) - \frac{|A|}{|C|} \left[ \langle \phi^+ | X | \phi^+ \rangle | \phi^+ \rangle \langle \phi^+ | \right. \right. \\ & \quad \left. \left. + \text{Tr}(\Pi^- X) \frac{1}{|C|^2 - 1} (|C|^2 D_K(\tau_C)_A \otimes \tau_{A'} - \phi_{AA'}^+) \right] \right\|_{\diamond} \\ & \leq 2\sqrt{2}\varepsilon|A| \left( 2\sqrt{|A||C|} + 1 \right). \end{aligned} \quad (\text{B.1})$$

*Proof.* It follows directly from the fact that  $(\tau_K, E, D)$  is a QES together with [Lemma A.1](#) that

$$\mathcal{E}(\phi_{CC'}^+) = \frac{|A|}{|C|} \phi_{AA'}^+. \quad (\text{B.2})$$

Let  $\Lambda_{C \rightarrow C\tilde{B}_1}^{(i)}$ ,  $i = 0, 1$  be two attack maps such that  $\eta_{\Lambda^{(i)}} |\phi^+\rangle = 0$  for  $i = 0, 1$  and define

$$\Lambda_{C \rightarrow C\tilde{B}_1\tilde{B}_2} = \frac{1}{2} \sum_{i=0,1} |i\rangle \langle i|_{\tilde{B}_2} \otimes \Lambda^{(i)}.$$

The the  $\varepsilon$ -NM property implies

$$I(AA' : \tilde{B}_1\tilde{B}_2)_{\eta_{\tilde{\Lambda}}} \leq \varepsilon,$$

and therefore, using Pinsker's inequality,

$$\begin{aligned} & \left\| \frac{1}{2} \sum_{i=0,1} |i\rangle \langle i|_{\tilde{B}} \otimes (\eta_{\tilde{\Lambda}^{(i)}})_{CC'\tilde{B}_1} \right. \\ & \quad \left. - \frac{1}{4} \left( \sum_{i=0,1} |i\rangle \langle i|_{\tilde{B}} \otimes (\eta_{\tilde{\Lambda}^{(i)}})_{\tilde{B}_1} \right) \otimes \left( \sum_{i=0,1} (\eta_{\tilde{\Lambda}^{(i)}})_{CC'} \right) \right\|_1 \leq \sqrt{2\varepsilon}. \end{aligned} \quad (\text{B.3})$$

Observe that

$$\begin{aligned} \eta_{\tilde{\Lambda}} &= \frac{1}{|K|} \sum_k D_k \circ \Lambda \circ E_k(\phi_{AA'}^+) \\ &= \frac{|C|}{|A|} \frac{1}{|K|} \sum_k (D_k \otimes E_k^T) \circ \Lambda(\phi_{CC'}^+) \\ &= \frac{|C|}{|A|} \mathcal{E} \circ \Lambda(\phi_{CC'}^+). \end{aligned} \quad (\text{B.4})$$

Setting  $(\eta_{A^{(0)}})_{CC'\bar{B}_1} = \tau_{CC'}^- \otimes (\eta_{A^{(1)}})_{\bar{B}_1}$ , we get

$$\begin{aligned}\eta_{A^{(0)}} &= \frac{|C|}{|A|} \mathcal{E}(\tau^-) \otimes (\eta_{A^{(1)}})_{\bar{B}_1} \\ &= \frac{|C|}{|A|} \frac{1}{|C|^2 - 1} (|C|^2 \mathcal{E}(\tau_{CC'}) - \mathcal{E}(\phi_{CC'}^+)) \otimes (\eta_{A^{(1)}})_{\bar{B}_1} \\ &= \frac{1}{|C|^2 - 1} (|C|^2 D_K(\tau_C) \otimes \tau_A - \phi_{AA'}^+) \otimes (\eta_{A^{(1)}})_{\bar{B}_1}.\end{aligned}\quad (\text{B.5})$$

and therefore

$$\begin{aligned}&\left\| \frac{1}{|C|^2 - 1} (|C|^2 D_K(\tau_C) \otimes \tau_A - \phi_{AA'}^+) \otimes (\eta_{A^{(1)}})_{\bar{B}_1} - \frac{|C|}{|A|} \mathcal{E}((\eta_{A^{(1)}})_{CC'\bar{B}_1}) \right\|_1 \\ &\leq 2\sqrt{2\varepsilon}\end{aligned}\quad (\text{B.6})$$

for all  $A^{(1)}$ . For any state  $\varrho_{CC'\bar{B}_1}$  with  $\varrho_{CC'\bar{B}}|\phi^+\rangle_{CC'} = 0$ , we define the state

$$\begin{aligned}\varrho'_{CC'\bar{B}_1\bar{B}_2} &\frac{1}{C} \left( |0\rangle\langle 0|_{\bar{B}_2} \otimes \varrho_{CC'\bar{B}_1} \right. \\ &\left. + |1\rangle\langle 1|_{\bar{B}_2} \otimes [((\mathbb{1}_C - \varrho_C) \otimes V_{C'}) \phi^+ ((\mathbb{1}_C - \varrho_C) \otimes V_{C'})] \otimes \varrho_{\bar{B}_2} \right).\end{aligned}\quad (\text{B.7})$$

Here,  $V$  is a unitary such that  $\text{Tr}(\mathbb{1}_C - \varrho_C)V_C^T = 0$ . It is easy to see that such a unitary always exists, the existence is equivalent to the fact that any  $|C|$ -tuple of real numbers is the ordered list of side lengths of a polygon in the complex plain. Note that  $\varrho'_{CC'\bar{B}_1\bar{B}_2}|\phi^+\rangle_{CC'} = 0$ , and  $\varrho'_C = \tau_{C'}$ . Together with the triangle inequality, equation (B.6) implies therefore that

$$\begin{aligned}&\frac{1}{|C|} \left\| \frac{|C|}{|A|} \mathcal{E}(\varrho) - \frac{1}{|C|^2 - 1} (|C|^2 D_K(\tau_C) \otimes \tau_A - \phi_{AA'}^+) \otimes \varrho_{\bar{B}_1} \right\|_1 \\ &+ \left\| \frac{|C|}{|A|} \mathcal{E} [((\mathbb{1}_C - \varrho_C) \otimes V_{C'}) \phi^+ ((\mathbb{1}_C - \varrho_C) \otimes V_{C'})] \right. \\ &\quad \left. - \frac{|C| - 1}{|C|} \frac{1}{|C|^2 - 1} (|C|^2 D_K(\tau_C) \otimes \tau_A - \phi_{AA'}^+) \right\|_1 \leq 2\sqrt{2\varepsilon},\end{aligned}$$

i.e. in particular

$$\left\| \frac{|C|}{|A|} \mathcal{E}(\varrho) - \frac{1}{|C|^2 - 1} (|C|^2 D_K(\tau_C) \otimes \tau_A - \phi_{AA'}^+) \otimes \varrho_{\bar{B}_1} \right\|_1 \leq 2\sqrt{2\varepsilon}|C|.$$

As  $\varrho$  was arbitrary we have proven that

$$\left\| \frac{|C|}{|A|} \mathcal{E} - \left\langle \frac{1}{|C|^2 - 1} (|C|^2 D_K(\tau_C) \otimes \tau_A - \phi_{AA'}^+) \right\rangle \right\|_{\diamond} \leq 2\sqrt{2\varepsilon}|C|. \quad (\text{B.8})$$

The only fact that is left to show is, that  $\|\mathcal{E}(|\phi^+\rangle\langle v|)\|_1$  is small for all normalized  $|v\rangle$  such that  $\langle \phi^+ | v \rangle = 0$ . To this end, observe that  $\text{Tr}_A \circ \mathcal{E}(\sigma_C \otimes (\cdot)_{C'}) = E_K^T$

for all quantum states  $\sigma_C$ . Let  $\varrho_C$  be any quantum state that does not have full rank, note that such states span all of  $\mathcal{B}(\mathcal{H}_C)$ , and for hermitian operators there exists a decomposition into such operators that saturates the triangle inequality. Taking a quantum state  $\sigma_C$  such that  $\langle \phi^+ | \varrho \otimes \sigma | \phi^+ \rangle = \frac{1}{|C|} \text{Tr} \varrho_C \sigma_C^T = 0$  (the first equality is the mirror lemma A.1), we have

$$\left\| \mathcal{E}(\varrho \otimes \sigma) - \frac{|A|}{|C|} \frac{1}{|C|^2 - 1} (|C|^2 D_K(\tau_C) \otimes \tau_A - \phi_{AA'}^+) \right\|_1 \leq 2\sqrt{2\varepsilon}|A|$$

according to what we have already proven. Using inequality (B.8) we arrive at

$$\left\| E_K^\dagger(X) - \frac{|A|}{|C|} \tau_A \text{Tr}(X) \right\|_1 \leq 2\sqrt{2\varepsilon}|A| \|X\|_1 \quad (\text{B.9})$$

For Hermitian matrices  $X$  and therefore

$$\left\| E_K^\dagger(X) - \frac{|A|}{|C|} \tau_A \text{Tr}(X) \right\|_1 \leq 4\sqrt{2\varepsilon}|A| \|X\|_1 \quad (\text{B.10})$$

For arbitrary  $X$ . We can write  $|v\rangle_{CC'} = X_C |\phi^+\rangle_{CC'}$  for some traceless matrix  $X_C$ . Now we calculate

$$\begin{aligned} \|\mathcal{E}(|\phi^+\rangle\langle v|_{CC'})\|_1 &= \left\| \frac{|A|}{|C|} |\phi^+\rangle\langle \phi^+|_{AA'} \left( E_K^\dagger(X^\dagger) \right)_A \right\|_1 \\ &= \frac{|A|}{|C|} \left\| \left( E_K^\dagger(X) \right)_A |\phi^+\rangle_{AA'} \right\|_2 \\ &= \frac{\sqrt{|A|}}{|C|} \|E_K^\dagger(X)\|_2 \\ &\leq \frac{\sqrt{|A|}}{|C|} \|E_K^\dagger(X)\|_1 \\ &\leq \frac{|A|^{3/2}}{|C|} 4\sqrt{2\varepsilon} \|X\|_1 \\ &\leq 4\sqrt{2\varepsilon}|A|^{3/2}. \end{aligned} \quad (\text{B.11})$$

The first equation is Lemma B.1, the second and third equations are easily verified, the first inequality is a standard norm inequality, the second inequality is Equation (B.10), and the last inequality follows from the normalization of  $|v\rangle$ . By the Schmidt decomposition, we get a stabilized version of this inequality,

$$\|\mathcal{E}(|\phi^+\rangle_{CC'} |\alpha\rangle_{\bar{B}_1} \langle v|_{CC' \bar{B}_1})\|_1 \leq 2\sqrt{2\varepsilon}|A|^{3/2}, \quad (\text{B.12})$$

for all  $|\alpha\rangle_{\tilde{B}_1}$  and all  $|v\rangle_{CC'\tilde{B}}$  such that  $\langle\phi^+|v\rangle = 0$  Combining everything we arrive at

$$\begin{aligned} & \left\| \mathcal{E}(X) - \frac{|A|}{|C|} \left[ \langle\phi^+|X|\phi^+\rangle|\phi^+\rangle\langle\phi^+| \right. \right. \\ & \quad \left. \left. + \text{Tr}(\Pi^-X) \frac{1}{|C|^2-1} (|C|^2 D_K(\tau_C)_A \otimes \tau_{A'} - \phi_{AA'}^+) \right] \right\|_{\diamond} \\ & \leq 2\sqrt{2\varepsilon}|A| \left( 4\sqrt{|A|} + 1 \right). \end{aligned} \quad (\text{B.13})$$

□

We are now ready to prove the characterization theorem [Theorem 4.4](#) in the  $\varepsilon$ -approximate setting (including the exact case, [Theorem 3.7](#) by setting  $\varepsilon = 0$ .)

**Theorem B.3 (Precise version of [Theorem 4.4](#))** *Let  $\Pi = (\tau, E, D)$  be a QES.*

1. *If  $\Pi$  is  $\varepsilon$ -NM, then any attack map  $\Lambda_{CB \rightarrow C\tilde{B}}$  results in an effective map  $\tilde{A}_{AB \rightarrow A\tilde{B}}$  fulfilling*

$$\left\| \tilde{A}_{AB \rightarrow A\tilde{B}} - \tilde{A}_{AB \rightarrow A\tilde{B}}^{\text{exact}} \right\|_{\diamond} \leq 2\sqrt{2\varepsilon}|A|^4|C| \left( 4\sqrt{|A|} + 1 \right), \quad (\text{B.14})$$

where

$$\tilde{A}_{AB \rightarrow A\tilde{B}}^{\text{exact}} = \text{id}_A \otimes \Lambda'_{B \rightarrow \tilde{B}} + \frac{1}{|C|^2-1} (|C|^2 \langle D_K(\tau) \rangle - \text{id})_A \otimes \Lambda''_{B \rightarrow \tilde{B}},$$

with  $\Lambda' = \text{Tr}_{CC'}[\phi_{CC'}^+, \Lambda(\phi_{CC'}^+ \otimes (\cdot))]$  and  $\Lambda'' = \text{Tr}_{CC'}[\Pi_{CC'}^-, \Lambda(\phi_{CC'}^+ \otimes (\cdot))]$ .

2. *Conversely, if for a scheme all effective maps fulfil Equation (B.14) with the right hand side replaced by  $\varepsilon$ , then it is  $5\varepsilon(\log(|A|) + r) + 3h(\varepsilon)$ -NM, where  $r$  is a bound on the size of the honest user's side information.*

*Proof.* We start with 1. We want to bound the diamond norm distance between the effective map  $\tilde{A}$  resulting from an attack  $\Lambda$  and the idealized effective map  $\tilde{A}^{\text{exact}}$ . Let

$$|\psi\rangle_{AA'BB'} = \sum_{i=0}^{|A|^2-1} \sqrt{p_i} |\alpha_i\rangle_{AA'} \otimes |\beta_i\rangle_{BB'}$$

be an arbitrary pure state given in its Schmidt decomposition across the bipartition  $AA'$  vs.  $BB'$ . We can write  $|\alpha_i\rangle_{AA'} = X_{A'}^{(i)} |\phi^+\rangle$  for some matrices  $X^{(i)}$  satisfying  $\|X^{(i)}\|_{\infty} \leq |A|$ . We calculate the action of  $\tilde{A}$  on  $|\alpha_i\rangle\langle\alpha_j|_{AA'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}$ ,

$$\begin{aligned} & \tilde{A}_{AB \rightarrow A\tilde{B}}^{\text{exact}} (|\alpha_i\rangle\langle\alpha_j|_{AA'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) = X_{A'}^{(i)} \left( |\phi^+\rangle\langle\phi^+|_{AA'} \otimes \Lambda'_{B \rightarrow \tilde{B}} (|\beta_i\rangle\langle\beta_j|_{BB'}) \right) \\ & + \frac{1}{|C|^2-1} (|C|^2 D_K(\tau)_A \otimes \tau_{A'} - |\phi^+\rangle\langle\phi^+|_{AA'}) \otimes \Lambda''_{B \rightarrow \tilde{B}} (|\beta_i\rangle\langle\beta_j|_{BB'}) X_{A'}^{(j)}. \end{aligned} \quad (\text{B.15})$$

In a similar way we get

$$\begin{aligned}
& \tilde{\Lambda}_{AB \rightarrow A\tilde{B}}(|\alpha_i\rangle\langle\alpha_j|_{AA'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) \\
&= X_{A'}^{(i)} \tilde{\Lambda}_{AB \rightarrow A\tilde{B}}(|\phi^+\rangle\langle\phi^+|_{AA'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) X_{A'}^{(i)} \\
&= \frac{|C|}{|A|} X_{A'}^{(i)} \mathcal{E}_{CC' \rightarrow AA'} \circ \Lambda_{CB \rightarrow C\tilde{B}}(|\phi^+\rangle\langle\phi^+|_{CC'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) X_{A'}^{(i)}. \quad (\text{B.16})
\end{aligned}$$

Using [Lemma B.2](#) we bound

$$\begin{aligned}
& \left\| \left( \tilde{\Lambda}_{AB \rightarrow A\tilde{B}} - \tilde{\Lambda}_{AB \rightarrow A\tilde{B}}^{\text{exact}} \right) (|\alpha_i\rangle\langle\alpha_j|_{AA'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) \right\|_1 \\
&= \left\| X_{A'}^{(i)} \left( \tilde{\Lambda}_{AB \rightarrow A\tilde{B}} - \tilde{\Lambda}_{AB \rightarrow A\tilde{B}}^{\text{exact}} \right) (|\phi^+\rangle\langle\phi^+|_{AA'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) X_{A'}^{(i)} \right\|_1 \\
&\leq \left\| X_{A'}^{(i)} \right\|_\infty \left\| X_{A'}^{(j)} \right\|_\infty \left\| \left( \tilde{\Lambda}_{AB \rightarrow A\tilde{B}} - \tilde{\Lambda}_{AB \rightarrow A\tilde{B}}^{\text{exact}} \right) (|\phi^+\rangle\langle\phi^+|_{AA'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) \right\|_1 \\
&= \left\| X_{A'}^{(i)} \right\|_\infty \left\| X_{A'}^{(j)} \right\|_\infty \left\| \frac{|C|}{|A|} \mathcal{E}_{CC' \rightarrow AA'} \circ \Lambda_{CB \rightarrow C\tilde{B}}(|\phi^+\rangle\langle\phi^+|_{CC'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) \right. \\
&\quad \left. - \tilde{\Lambda}_{AB \rightarrow A\tilde{B}}^{\text{exact}}(|\phi^+\rangle\langle\phi^+|_{AA'} \otimes |\beta_i\rangle\langle\beta_j|_{BB'}) \right\|_1 \\
&\leq 2\sqrt{2\varepsilon}|A|^2|C| \left( 4\sqrt{|A|} + 1 \right). \quad (\text{B.17})
\end{aligned}$$

The inequalities result from applying Hölder's inequality twice, and [Lemma B.2](#), respectively. Using the triangle inequality we get

$$\begin{aligned}
& \left\| \left( \tilde{\Lambda}_{AB \rightarrow A\tilde{B}} - \tilde{\Lambda}_{AB \rightarrow A\tilde{B}}^{\text{exact}} \right) (|\psi\rangle\langle\psi|_{AA'BB'}) \right\|_1 \\
&\leq 2\sqrt{2\varepsilon}|A|^2|C| \left( 4\sqrt{|A|} + 1 \right) \sum_{i,j=0}^{|A|^2-1} \sqrt{p_i p_j} \\
&\leq 2\sqrt{2\varepsilon}|A|^4|C| \left( 4\sqrt{|A|} + 1 \right). \quad (\text{B.18})
\end{aligned}$$

As  $|\psi\rangle$  was arbitrary, we have proven

$$\left\| \tilde{\Lambda}_{AB \rightarrow A\tilde{B}} - \tilde{\Lambda}_{AB \rightarrow A\tilde{B}}^{\text{exact}} \right\|_\diamond \leq 2\sqrt{2\varepsilon}|A|^4|C| \left( 4\sqrt{|A|} + 1 \right). \quad (\text{B.19})$$

Now let us prove 2. Let  $\Lambda_{CB \rightarrow C\tilde{B}}$  again be an arbitrary attack map, and assume that the resulting effective map is  $\varepsilon$ -close to  $\tilde{\Lambda}_{AB \rightarrow A\tilde{B}}^{\text{exact}}$ . Observe that  $p^\varepsilon(\Lambda, \varrho) = \text{Tr} \Lambda'(\varrho_B)$ .

By the Alicki-Fannes inequality [4] and [Lemma A.2](#), this implies

$$I(\Lambda R : \tilde{B})_{\tilde{\Lambda}(\varrho)} \leq I(\Lambda R : B)_\varrho + h(p^\varepsilon(\Lambda, \varrho)) + 5\varepsilon \log(|A||R|) + 3h(\varepsilon) \quad (\text{B.20})$$

with the help of [Lemma A.2](#).  $\square$