

# Actively Secure Two-Party Evaluation of any Quantum Operation

Frédéric Dupuis<sup>1\*</sup> Jesper Buus Nielsen<sup>2\*\*</sup> Louis Salvail<sup>3\*\*\*</sup>

<sup>1</sup> Institute for Theoretical Physics, ETH Zurich, Switzerland  
dupuis@phys.ethz.ch

<sup>2</sup> Department of Computer Science, Aarhus University, Denmark  
jbn@cs.au.dk

<sup>3</sup> Université de Montréal (DIRO), QC, Canada  
salvail@iro.umontreal.ca

**Abstract.** We provide the first two-party protocol allowing Alice and Bob to evaluate privately even against active adversaries any completely positive, trace-preserving map  $\mathcal{F} \in \mathcal{L}(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}) \rightarrow \mathcal{L}(\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}})$ , given as a quantum circuit, upon their joint quantum input state  $\rho_{\text{in}} \in \mathcal{D}(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}})$ . Our protocol leaks no more to any active adversary than an ideal functionality for  $\mathcal{F}$  provided Alice and Bob have the cryptographic resources for active secure two-party classical computation. Our protocol is constructed from the protocol for the same task secure against specious adversaries presented in [DNS10].

## 1 Introduction

We provide the first active-secure two-party protocol for computing on quantum data. We look at a model where Alice and Bob hold an input  $\rho_{\text{in}}$  on registers  $\mathcal{A}_{\text{in}}$  and  $\mathcal{B}_{\text{in}}$ , where Alice holds register  $\mathcal{A}_{\text{in}}$  and Bob holds  $\mathcal{B}_{\text{in}}$ . They agree on a completely positive, trace-preserving (CPTP) map  $\mathcal{F}$  from registers  $\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}$  to registers  $\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}}$ , and they want to compute  $\rho_{\text{out}} = \mathcal{F}(\rho_{\text{in}})$  such that, at the end of the protocol, Alice is in possession of  $\mathcal{A}_{\text{out}}$  and Bob is in possession of  $\mathcal{B}_{\text{out}}$ . They want to do this in an actively secure manner. Our notion of active security is phrased via simulation, but intuitively it simply guarantees that any cheating Alice, even an infinitely powerful Alice, which might deviate from the protocol, can only affect the output of the protocol by replacing her own input and that she will at any point during the execution of the protocol only hold information which can be computed (efficiently) from either  $\rho_{\text{in}}^{\mathcal{A}}$  or  $\rho_{\text{out}}^{\mathcal{A}}$ . The equivalent condition should hold for Bob.

A simple example of such an  $\mathcal{F}$  is the quantum swap, where  $\mathcal{A}_{\text{in}} = \mathcal{A}_{\text{out}} = \mathcal{A}$ ,  $\mathcal{B}_{\text{in}} = \mathcal{B}_{\text{out}} = \mathcal{B}$ , and  $\rho_{\text{out}}^{\mathcal{A}} = \rho_{\text{in}}^{\mathcal{B}}$  and  $\rho_{\text{out}}^{\mathcal{B}} = \rho_{\text{in}}^{\mathcal{A}}$ . Securely implementing this unitary basically means to do an atomic swap of quantum states, which was shown impossible in [DNS10] even against a restricted class of adversaries called *specious*. Extra assumptions are therefore needed to get unconditional security. Our way out is to look at a model where the two parties have access to an ideal functionality which allows them to securely do any *classical* computation on any *classical*

---

\* Supported by Canada's NSERC Postdoctoral Fellowship Program and by the Swiss National Science Foundation via the National Center of Competence QSIT.

\*\* Partially supported by an European Research Council, Starting Grant, number 279447 and the Danish National Research Foundation and the National Science Foundation of China (under the grant 61061130540) for the Sino-Danish Center for the Theory of Interactive Computation.

\*\*\* Supported by Canada's NSERC discovery grant, FREQUENCY(NSERC), the QuantumWorks networks(NSERC), and INTRIQ(FQRNT).

data held jointly by the two parties. In this model we give an unconditionally secure protocol with active security. This is the first such protocol.

Formally, we use the notationally simpler model of [DNS10], but it is easy to see that security in this model implies security in the stand-alone model of [HSS11], as long as the simulator is poly-time. The stand-alone model of [HSS11] allows, inside a secure quantum protocol, to replace a classical ideal functionality by a classical protocol which securely implements that ideal functionality against poly-time quantum adversaries. The result is a protocol secure against poly-time quantum adversaries. This, in particular, implies sequential security, i.e., if a protocol is secure, it remains secure when run in sequence with other secure protocols. Since the secure evaluation of any classical function with security against poly-time quantum adversaries can be done under the assumption that learning with errors is hard [Reg05,HSS11] and under the more general assumption that mixed commitment schemes exist [LN11], our poly-time simulator provides an active-secure two-party plain-model protocol for computing on quantum data with security against any poly-time quantum adversaries. This is the first such protocol.

### 1.1 Overview of our construction

We reuse many ideas from the protocol provided in [DNS10], which gives a two-party protocol for computing on quantum data securely against so-called specious adversaries. The protocol therein is unconditionally secure given an ideal functionality for classical computation. Specious adversaries are a quantum version of the classical notion of passive adversaries. Technically, a specious adversary is an adversary which is allowed to deviate from the protocol, except that at any step of the protocol it should be able to reconstruct the honest state of the protocol from its current state. This basically allows it to purify itself and not much else.

In the protocol in [DNS10], all wires are encrypted using a Pauli encryption: a qubit  $|v\rangle$  is represented as  $|V\rangle = X^x Z^z |v\rangle$ , where the two uniform key bits  $x$  and  $z$  are secret-shared between Alice and Bob. For example, to secret-share  $x$ , Alice will hold a uniformly random bit  $x_A$  and Bob will hold a uniformly random bit  $x_B$  such that  $x = x_A \oplus x_B$ . All wires are independently encrypted like this, which ensures that intermediate states are perfectly hidden from both parties. Computation is then done “through the encryption”. The CPTP map  $\mathcal{F}$  is described by a quantum circuit made out of the universal set of gates  $\mathcal{UG}$  consisting of gates X, Y, Z, CNOT,  $H = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$ ,  $P = \begin{pmatrix} 1 & 0 \\ 0 & i \end{pmatrix}$ , and  $R = \begin{pmatrix} 1 & 0 \\ 0 & e^{i\pi/4} \end{pmatrix}$ , together with a set of ancilla wires initialized in state  $|0\rangle$ , and a set of output wires that are discarded (or, in our case, that remain encrypted forever). The protocol evaluates each gate of  $\mathcal{F}$  while preserving privacy. Handling the Pauli gates X, Y and Z is easy, as they commute or anti-commute with the encryption operators. As an example, assume that the parties want to apply an X gate to a qubit  $|v\rangle$ , i.e., compute an encryption of  $|v'\rangle = X|v\rangle$ . Since up to an overall phase factor  $XZ = ZX$ , it follows that  $X|V\rangle = XX^x Z^z |v\rangle = X^x Z^z X|v\rangle = X^x Z^z |v'\rangle$ . So, the evaluation is simply performed on the encrypted qubit  $|V\rangle$  and the key bits  $x$  and  $z$  are maintained. For the remaining Clifford gates CNOT, H and P other “commutation” rules with X and Z are used. For H, it is used that  $HX = ZH$  and  $HZ = XH$ , so H is simply applied to the encrypted qubit, and then the keys  $x$  and  $z$  are swapped: Alice sets  $x'_A = z_A$  and  $z'_A = x_A$  and similarly for Bob. This leaves the non-Clifford gate R, where the relation  $RX^x Z^z = P^x X^x Z^z R$  almost does the job, except that it leaves an extra  $P^x$ . Getting rid of this requires quantum communication and a classical secure two-party computation, which computes how the parties should update their key bits. After such a gate-by-gate computation of  $U$  on the encrypted qubits, the

shares of the keys needed for learning the states on ones own output wires are swapped with the other party in one atomic step, using the ideal functionality for classical secure computation.

The protocol in [DNS10] is actually secure against an active adversary up to the final swap of key bits, as the encryptions of the wires are guaranteed to be perfect, independent of the behavior of the other party, as the parties pick their own shares of the sharings of  $x$  and  $z$ . This means that no party can get any information on any intermediary states, no matter how it deviates. It can, however, easily force the computation to be incorrect, by applying gates to the encrypted states, so the full protocol is not active secure. We note, however, that [DNS10] is secure against what we could call *ultimately specious adversaries*: Adversaries promising to attack in such a way that both parties always be able to reconstruct the correct state at the *end* of the protocol, but that can otherwise behave as they want. This follows from the active security in the middle and a theorem in [DNS10] which says that any attack which always allows both parties to obtain the correct output can be simulated given just the output of the corrupted party—basically, there is no way to learn extra information without sometimes irrevocably destroying the state of the other party.

In this paper, we use this observation in a protocol proceeding along the same lines as [DNS10] but where we force the adversary to be ultimately specious. This is done by not only encrypting the wires, but by unconditionally authenticating them. In addition, we commit the parties to their key bits, to allow the recipient to verify the key bits swapped at the end. Since an unconditional quantum authentication code is also an unconditionally secure encryption, we get a protocol with at least the security of [DNS10], but with the added property that an adversary who deviates from the protocol will be caught. More technically, if all checks of the authentication code succeed, then the authenticated values collapse to the correct values. This forces the adversary to either be detected or be ultimately specious. Since we do all the checks of the authentication codes *before* any key bits are revealed, the case where the adversary is detected can be simulated by simply asking the ideal world to abort too. The case where the adversary is ultimately specious is simulated similarly to [DNS10].

The main technical challenge is then to devise an authentication code with these two properties:

1. It allows to perform computation “through the authentication”.
2. It allows to check the authentication code without revealing what is authenticated.

The first property is important for hiding intermediate values during the computation. The second property is important when we force the adversary to either be detected or ultimately specious: when he is detected, he should learn nothing on the incorrect outputs.

We devise an authentication code with these properties based on the Clifford-based quantum authentication code proposed in [ABE08]. It authenticates a quantum message using a random unitary implementable using gates X, Y, Z, CNOT, H and P. The authentication works as follows. In order to satisfy the second property, take a qubit  $|v\rangle$  on some wire. Alice prepends  $n$  new wires, in the all-zero state  $|0^n\rangle$ . Then she applies a uniformly random  $(n + 1)$ -bit Clifford operator  $A$  to the  $n + 1$  wires. She then sends the state to Bob, who appends  $n$  more wires in the all-zero state and applies a uniformly random  $(2n + 1)$ -bit Clifford operator  $B$  to the  $2n + 1$  wires. We can write this as  $|V\rangle = B(|0^n\rangle \otimes A(|v\rangle|0^n\rangle))$ . The key is  $(A, B)$ . This authentication can be checked in two different ways. Either, apply  $B^\dagger$  to the authenticated state and check that the first  $n$  wires are all zero. Or, apply  $B^\dagger$  and then apply  $A^\dagger$  to the last  $n + 1$  wires and check that the last  $n$  wires are all zero. One way is used for Alice to check that Bob did not change the authenticated value. The

other is used by Bob to check Alice. In order to perform these without leakage, we use a 2-party classical ideal functionality as described below.

In our scheme, Alice will hold  $A$  as her share of the key and Bob will hold  $B$  as his key share. They are committed to their share by being committed to a poly-size classical description of the operator applied. We sketch why the scheme has the two necessary properties.

1. Since the authentication is performed using only Clifford gates, an approach as in [DNS10] will allow to fairly easily apply Clifford gates “through the authentication”. Since the Clifford unitaries form a group, the operation consisting of decrypting a qubit, applying a Clifford gate to it and reauthenticating it using a different key is also Clifford unitary. Hence, we can apply a Clifford gate to a swaddled qubit simply by changing the authentication keys in the appropriate manner. More precisely, to execute the Clifford gate  $G \in \mathcal{C}_1$ , Alice updates her key to  $A(G^\dagger \otimes \mathbb{1}_n)$ , and Alice and Bob use a TPC to update Alice’s commitment to her key. Executing CNOT gates is similar, but involves two swaddlings. The R-gate requires new techniques reminiscent of the fault tolerant implementation of it [Sho96,Got97,Got10]. The details appear within, but we basically reduce it to securely producing a state  $|0\rangle$ , a magic state, a measurement in the computational basis, and applying secure Clifford gates together with a carefully chosen secure classical computation which tells the parties how to update their keys.
2. If Bob is to check an authenticated qubit, we simply give him  $|V\rangle$  and he applies  $B^\dagger$  and measures the first  $n$  wires, rejecting if they are not all zero. After that he re-applies  $B$  to recover the authentication  $|V\rangle$ . If Alice is to perform the check, we perform a secure classical computation where the inputs are the committed descriptions of  $A$  and  $B$  which effectively swaps the inner and outer authentications (see below for details). Alice then holds the outermost authentication and can easily test its integrity.

The final state of the computation is obtained after each party reveals the authentication keys needed by the other party to open its output wires.

## 2 Preliminaries

We refer to [DNS10] for a more complete introduction to the notation used throughout. In the following, various Hilbert spaces will be seen as collections of *wires*, represented in lower-case typewriter font (such as  $w, x, y$ ); each wire represents a qubit involved in a computation. We will denote by  $w \in \mathcal{A}$  the fact that  $w$  is part of the system  $\mathcal{A}$ . In the following, registers like  $\mathcal{A}$  are considered as being composed by a set of wires or qubit registers. Henceforth, we shall denote a negligible function in  $n$  by  $\text{negl}(n)$ .

The set of hermitian operators on  $\mathcal{A}$  is denoted by  $\text{Herm}(\mathcal{A})$ , the set of positive semi-definite operators on  $\mathcal{A}$  is denoted by  $\text{Pos}(\mathcal{A})$  and the one with trace 1 are denoted by  $\text{D}(\mathcal{A})$ ;  $\text{D}(\mathcal{A})$  is the set of all possible quantum states for register  $\mathcal{A}$ . An operator  $A \in \text{L}(\mathcal{A}, \mathcal{B})$  is called a *linear isometry* if  $A^\dagger A = \mathbb{1}_{\mathcal{A}}$ . The set of unitary operators (i.e., linear isometries with  $\mathcal{B} = \mathcal{A}$ ) acting in  $\mathcal{A}$  is denoted by  $\text{U}(\mathcal{A})$ . The identity operator in  $\mathcal{A}$  is denoted  $\mathbb{1}_{\mathcal{A}}$  while the completely mixed state in  $\text{D}(\mathcal{A})$  is denoted by  $\mathbb{I}_{\mathcal{A}}$ . For any positive integer  $N > 0$ ,  $\mathbb{1}_N$  and  $\mathbb{I}_N$  denote the identity operator respectively the completely mixed state in the  $N$ -dimensional Hilbert space  $\mathcal{H}_N$ . When the context requires, a pure state  $|\psi\rangle \in \mathcal{A} \otimes \mathcal{B}$  will be written  $|\psi\rangle^{\mathcal{A}\mathcal{B}}$  to make explicit the registers in which it is stored. For Hilbert spaces  $\mathcal{A}$  and  $\mathcal{B}$ , we write  $\mathcal{A} \approx \mathcal{B}$  whenever  $\dim \mathcal{A} = \dim \mathcal{B}$ .

A linear mapping  $\Phi : L(\mathcal{A}) \mapsto L(\mathcal{B})$  is called a *super-operator*.  $\Phi$  is said to be *positive* if  $\Phi(A) \in \text{Pos}(\mathcal{B})$  for all  $A \in \text{Pos}(\mathcal{A})$ . The super-operator  $\Phi$  is said to be *completely positive* if  $\Phi \otimes \mathbb{1}_{L(\mathcal{Z})}$  is positive for every choice of the Hilbert space  $\mathcal{Z}$ . A super-operator  $\Phi$  can be physically realized or is *admissible* if it is completely positive and preserves the trace:  $\text{tr}(\Phi(A)) = \text{tr}(A)$  for all  $A \in L(\mathcal{A})$ . We call such a super-operator a *quantum operation* or a *CPTP map*. Any quantum operation  $\Phi : L(\mathcal{A}) \mapsto L(\mathcal{B})$  can be written in its Kraus form  $\{E_j\}_{j=1}^{\dim(\mathcal{A}) \cdot \dim(\mathcal{B})}$  where  $E_j \in L(\mathcal{A}, \mathcal{B})$  satisfies  $\sum_j E_j^\dagger E_j = \mathbb{1}_{\mathcal{A}}$  and  $\Phi(\rho) = \sum_j E_j \rho E_j^\dagger$ , for any  $\rho \in \text{Pos}(\mathcal{A})$ . Another, yet equivalent, way to represent any quantum operation is through a linear isometry  $W \in L(\mathcal{A}, \mathcal{B} \otimes \mathcal{Z})$  such that  $\Phi(\rho) = \text{tr}_{\mathcal{Z}}(W \rho W^\dagger)$ , for some extra space  $\mathcal{Z}$ . Any such isometry  $W$  can be implemented by a physical process as long as resources to implement the space  $\mathcal{Z}$  are available. This is just a unitary transform in  $U(\mathcal{A} \otimes \mathcal{Z})$  where the system in  $\mathcal{Z}$  is initially in known state  $|0\rangle^{\mathcal{Z}}$ . If  $\mathcal{E} : L(\mathcal{A}) \mapsto L(\mathcal{B})$  and  $\mathcal{F} : L(\mathcal{B}) \mapsto L(\mathcal{C})$  are two quantum operations (i.e. CPTP maps) then  $\mathcal{F} \circ \mathcal{E} : L(\mathcal{A}) \mapsto L(\mathcal{C})$  denotes the quantum operation that first applies  $\mathcal{E}$  then  $\mathcal{F}$ .

For  $\rho_0, \rho_1 \in D(\mathcal{A})$ , we denote by  $\Delta(\rho_0, \rho_1)$  the trace norm distance between  $\rho_0$  and  $\rho_1$ :  $\Delta(\rho_0, \rho_1) := \frac{1}{2} \|\rho_0 - \rho_1\|$ .

The Pauli group on one qubit is defined as  $P_1 := \{\pm \mathbb{1}, \pm i \mathbb{1}, \pm X, \pm i X, \pm Y, \pm i Y, \pm Z, \pm i Z\}$  while the  $n$ -qubit Pauli group is simply  $P_n := P_1^{\otimes n}$ . We denote by  $\mathcal{C}_n := \{U \in U(\mathcal{H}_{2^n}) : (\forall \sigma \in P_n)[U \sigma U^\dagger \in P_n]\}$ , the set of all Clifford operators acting on  $n$  qubits. The Clifford group is just the set of  $n$ -qubit unitary operators  $U$  with the property that for all Pauli operators  $\sigma \in P_n$ , we have  $\sigma' \in P_n$  such that  $U \sigma = \sigma' U$ .

## 2.1 Secure Two-Party Classical Computation Against Quantum Adversaries

Our protocol will use various *classical* two-party computations throughout its execution, each modeled as an ideal functionality. Recent work [LN11,HSS11] show that composable classical two-party computation protocols can be devised with security against quantum adversaries provided some classical computational assumptions hold against this class of adversaries. One example of such an assumption is *learning with errors is hard* [Reg05,HSS11]. The framework in [HSS11] allows us to replace the ideal functionalities with such secure protocols. Here we therefore focus on proving security given the ideal functionalities.

In the following, the ideal functionality for string commitment will be denoted by  $\text{id}_{\text{sc}}$ . It is defined as follows:

$$\begin{aligned} \text{id}_{\text{sc}}((id, s), \perp) &= (id, (id, \text{committed})) \text{ if } id \text{ is new } , \\ \text{id}_{\text{sc}}(\perp, (id, s)) &= ((id, \text{committed}), id) \text{ if } id \text{ is new } , \\ \text{id}_{\text{sc}}(s, s') &= (\perp, \perp) \text{ otherwise } . \end{aligned}$$

In order for Alice to commit on  $s \in \{0, 1\}^*$ , Alice and Bob call  $\text{id}_{\text{sc}}((id, s), \perp)$ , where  $id$  is an unused identifier chosen by Alice. In order for Bob to commit on  $s$ , Alice and Bob call  $\text{id}_{\text{sc}}(\perp, (id, s))$ , where  $id$  is chosen by Bob. The opening of a commitment is performed by calling the ideal functionality  $\text{id}_{\text{open}}$  defined as:

$$\begin{aligned} \text{id}_{\text{open}}(id, id) &= (s, s) \text{ if } \text{id}_{\text{sc}}((id, s), \perp) \text{ or } \text{id}_{\text{sc}}(\perp, (id, s)) \text{ was performed } , \\ \text{id}_{\text{sc}}(\cdot, \cdot) &= (\perp, \perp) \text{ otherwise } . \end{aligned}$$

Note that state has to be passed from  $\text{id}_{\text{sc}}(\cdot, \cdot)$  to  $\text{id}_{\text{OPEN}}(\cdot, \cdot)$  for the above descriptions to make sense. Our framework exactly allows that an ideal functionality from an earlier round passes its state to an ideal functionality in a later round. In the framework of [HSS11] they would be considered one ideal functionality, with a state. We prefer the above notation for brevity of later protocol descriptions.

Ideal functionalities computing functions applied to committed values can now be easily defined. Suppose that Alice and Bob want to let Alice learn a function  $f : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{0, 1\}^*$  upon two committed values,  $s$  and  $s'$ , where  $s$  is committed upon by Alice under  $\text{id}$  and  $s'$  is committed upon by Bob under  $\text{id}'$ . It suffices to define the ideal functionality  $\text{id}_f^*$  as

$$\begin{aligned} \text{id}_f^*((\text{id}, \text{id}'), (\text{id}, \text{id}')) &= (f(s, s'), \perp) \text{ if } \text{id}_{\text{OPEN}}(\text{id}, \text{id}) = (s, s) \wedge \text{id}_{\text{OPEN}}(\text{id}', \text{id}') = (s', s'), \\ \text{id}_f^*(\cdot, \cdot) &= (\perp, \perp) \text{ otherwise.} \end{aligned}$$

The same if it is Bob who is to learn the output, but with  $\text{id}_f^*((\text{id}, \text{id}'), (\text{id}, \text{id}')) = (\perp, f(s, s'))$ . The same construction can be used to implement any efficiently computable function  $f$  evaluated upon *any* number of committed values. The ideal functionality can also easily be extended to produce commitments to the outputs of  $f$ . It is this extended ideal functionality we use most often. Again,  $\text{id}_f^*(\cdot, \cdot)$  will have to share state with  $\text{id}_{\text{sc}}(\cdot, \cdot)$  and  $\text{id}_{\text{OPEN}}(\cdot, \cdot)$ , which is allowed in our framework.

Note that all the above ideal functionalities are defined such that at most one party has a non-trivial output (i.e., an output which is not known before the inputs are provided). We avoid using functions where both parties have a non-trivial output as an easy way to deal with the problem that fairness in classical secure two-party computation is provably impossible for most functionalities. It is not hard to see that it follows from known completeness results of quantum-secure classical two-party computation [LN11, HSS11] that the classical ideal functionalities specified above can be implemented with security against poly-time quantum adversaries. We skip the details of this, as our focus here is on using the classical ideal functionalities for constructing secure two-party protocols for computing on quantum data.

## 2.2 Clifford-Based Quantum Authentication

In order to detect misbehaviors of an active adversary, we will be evaluating a circuit upon authenticated quantum bits. We will be using a quantum authentication scheme (QAS) [BCG<sup>+</sup>02] based on Clifford operators introduced in [ABE08] as our main building block.

**Definition 2.1 (Quantum authentication scheme).** A quantum authentication scheme is a set of encryption and decryption superoperators  $\{(\mathcal{E}_k^{S \rightarrow C}, \mathcal{D}_k^{C \rightarrow SF}) : k \in \mathcal{K}\}$ , where  $\mathcal{K}$  is the set of possible keys,  $S$  is the input system,  $C$  is the ciphertext system, and  $F$  is a “flag” system that contains either  $|\text{acc}\rangle$  or  $|\text{rej}\rangle$ . A QAS is such that for all  $k \in \mathcal{K}$ ,  $(\mathcal{D}_k \circ \mathcal{E}_k)(\rho_S) = \rho_S \otimes |\text{acc}\rangle\langle \text{acc}|_F$ .

A QAS is secure if it satisfies the following:

**Definition 2.2 (Security of a QAS).** Let  $\mathcal{E}_k^{S \rightarrow C}$  and  $\mathcal{D}_k^{C \rightarrow SF}$  be the encoder and decoder corresponding to key  $k$ . Then, we say that the QAS  $(\mathcal{E}, \mathcal{D})$  is  $\varepsilon$ -secure if, for all attacks  $U_{CR}$ , there exists two CP maps  $\mathcal{U}_{R \rightarrow R}^{\text{acc}}$  and  $\mathcal{U}_{R \rightarrow R}^{\text{rej}}$  with  $\mathcal{U}^{\text{acc}} + \mathcal{U}^{\text{rej}} = \mathbb{1}$  such that for all inputs  $\psi_{SR}$ , we

have that for some fixed state  $\Omega_S$ :

$$\left\| \frac{1}{\#\mathcal{K}} \sum_{k \in \mathcal{K}} \mathcal{D}_k \left( U_{CR} \mathcal{E}_k(\psi_{SR}) U_{CR}^\dagger \right) - (\mathcal{W}^{\text{acc}}(\psi_{SR}) \otimes |\text{acc}\rangle\langle\text{acc}|_F + \mathcal{W}^{\text{rej}}(\psi_R) \otimes \Omega_S \otimes |\text{rej}\rangle\langle\text{rej}|_F) \right\|_1 \leq \varepsilon . \quad (1)$$

This definition can be shown to be equivalent to the existence of a simulator that interacts only with an ideal functionality in which Eve's only choice is whether or not to destroy the state and cause a rejection.

**Definition 2.3 (Clifford-based QAS[ABE08]).** *Let  $S$  be an  $s$ -qubit system,  $A$  be an  $n$ -qubit system, and let  $C = S \otimes A$ . Let  $\mathcal{K}$  index all Clifford unitaries  $C_k$  on  $s + n$  qubits. Then, the Clifford-based QAS is defined by the following encryption and decryption maps where  $P_{\text{acc}}^{SA} = \mathbb{1}_S \otimes |0^n\rangle\langle 0^n|_A$  and  $P_{\text{rej}}^{SA} = \mathbb{1}_{SA} - P_{\text{acc}}^{SA}$ :*

$$\begin{aligned} \mathcal{E}_k(\rho_S) &= C_k (\rho_S \otimes |0^n\rangle\langle 0^n|_A) C_k^\dagger , \\ \mathcal{D}_k(\sigma_{SA}) &= \text{tr}_A \left( P_{\text{acc}} C_k^\dagger \sigma_{SA} C_k P_{\text{acc}} \otimes |\text{acc}\rangle\langle\text{acc}|_F \right. \\ &\quad \left. + \text{tr} (P_{\text{rej}} C_k^\dagger \sigma_{SA} C_k) \pi_{SA} \otimes |\text{rej}\rangle\langle\text{rej}|_F \right) , \end{aligned}$$

where  $\pi_{SA}$  is an arbitrary fixed state that the decoder outputs when it rejects the authentication.

The following establishes the security of the QAS based on random Clifford operators. The proof of security is more or less the same as in [ABE08] and is provided in Appendix C for completeness:

**Theorem 2.4 (Security of Clifford-based QAS).** *The QAS defined above is  $\varepsilon(n)$ -secure for  $\varepsilon(n) = 6 \times 2^{-n}$ .*

It should be mentioned that picking a random Clifford operation acting upon  $\ell$  qubits requires to pick a uniformly random  $\text{poly}(\ell)$ -bit classical key  $k$  and the mapping between  $k$  and the corresponding Clifford operation can be performed efficiently [Got97, ABE08]. In other words, the key size of the Clifford-based QAS is polynomial in the number of qubits  $\ell = s + n$  used to authenticate an  $s$ -qubit quantum state.

### 2.3 Two-Party Quantum Protocols

We define two-party strategies in a similar way as in [DNS10, GW07], with some adaptations made for the fact that we are computing a CPTP map and not just a unitary and that we allow ideal functionalities of different rounds to share states (equivalent to considering one, stateful functionality). Two-party protocols for the evaluation of some CPTP map are particular cases of two-party strategies. Two-party strategies have access to some oracle in each round. An oracle is just a CPTP map acting on registers at both Alice and Bob. Oracles implement some functionalities like a communication channel or some more complex two-party functionalities.

An  $m$ -turn two-party strategy with oracle calls is defined by tuples of quantum operations  $\mathcal{A} := (\mathcal{A}_1, \dots, \mathcal{A}_{m+1})$ ,  $\mathcal{B} := (\mathcal{B}_1, \dots, \mathcal{B}_{m+1})$ , and  $\mathcal{O} := (\mathcal{O}_1, \dots, \mathcal{O}_m)$ . For  $i \in [1..m + 1]$ ,

operations  $\mathcal{A}_i \in L(\mathcal{A}_{i-1} \otimes \mathcal{A}_{i-1}^{\mathcal{O}, \text{out}}) \mapsto L(\mathcal{A}_i \otimes \mathcal{A}_i^{\mathcal{O}, \text{in}})$  and  $\mathcal{B}_i \in L(\mathcal{B}_{i-1} \otimes \mathcal{B}_{i-1}^{\mathcal{O}, \text{out}}) \mapsto L(\mathcal{B}_i \otimes \mathcal{B}_i^{\mathcal{O}, \text{in}})$  are the actions performed at turn  $i$  by Alice and Bob respectively. The operation  $\mathcal{O}_i : L(\mathcal{A}_i^{\mathcal{O}, \text{in}} \otimes \mathcal{O}_{i-1} \otimes \mathcal{B}_i^{\mathcal{O}, \text{in}}) \mapsto L(\mathcal{A}_i^{\mathcal{O}, \text{out}} \otimes \mathcal{O}_i \otimes \mathcal{B}_i^{\mathcal{O}, \text{out}})$  models the oracle provided to Alice and Bob at turn  $i$ , and  $\mathcal{O}_i$  a register used for passing state from the oracle of one round to the oracle at the next round. The oracle call at turn  $i \in [1..m]$  takes place right after  $\mathcal{A}_i$  and  $\mathcal{B}_i$  have been applied. In particular, these operations set the input registers  $\mathcal{A}_i^{\mathcal{O}, \text{in}}$  and  $\mathcal{B}_i^{\mathcal{O}, \text{in}}$  for the call to  $\mathcal{O}_i$ . The outputs are available to Alice and Bob in the next turn, in the output registers  $\mathcal{A}_i^{\mathcal{O}, \text{out}}$  and  $\mathcal{B}_i^{\mathcal{O}, \text{out}}$ . We make one exception to this general form, the last turn (that is turn  $m + 1$ ) of a strategy does not invoke any oracle, and is there simply to allow Alice and Bob to post-process the output of the last oracle.

Let  $\Pi = (\mathcal{A}, \mathcal{B}, \mathcal{O}, m)$  be an  $m$ -turn two-party strategy with oracle calls. The final state of the interaction between  $\mathcal{A}$  and  $\mathcal{B}$  upon joint input state  $\rho_{\text{in}} \in D(\mathcal{A}_0 \otimes \mathcal{B}_0 \otimes \mathcal{R})$ , where  $\mathcal{R}$  is a reference system with  $\dim \mathcal{R} = \dim \mathcal{A}_0 \dim \mathcal{B}_0$ , is denoted by

$$\begin{aligned} [\mathcal{A} \otimes \mathcal{B}]^{\mathcal{O}}(\rho_{\text{in}}) &:= (\mathcal{A}_{m+1} \otimes \mathcal{B}_{m+1} \otimes \mathbb{1}_{L(\mathcal{R} \otimes \mathcal{O}_{m+1})}) \\ &\quad (\mathbb{1}_{L(\mathcal{A}_m \otimes \mathcal{B}_m \otimes \mathcal{R})} \otimes \mathcal{O}_m)(\mathcal{A}_m \otimes \mathcal{B}_m \otimes \mathbb{1}_{L(\mathcal{R} \otimes \mathcal{O}_m)}) \\ &\quad \dots (\mathbb{1}_{L(\mathcal{A}_1 \otimes \mathcal{B}_1 \otimes \mathcal{R})} \otimes \mathcal{O}_1)(\mathcal{A}_1 \otimes \mathcal{B}_1 \otimes \mathbb{1}_{L(\mathcal{R})})(\rho_{\text{in}}) . \end{aligned}$$

A *communication oracle* from Alice to Bob is modeled by having  $\mathcal{A}_i^{\mathcal{O}, \text{in}} \approx \mathcal{B}_i^{\mathcal{O}, \text{out}}$  and letting  $\mathcal{O}_i$  move the state in  $\mathcal{A}_i^{\mathcal{O}, \text{in}}$  to  $\mathcal{B}_i^{\mathcal{O}, \text{out}}$ . A classical ideal functionality, as those described in Sect. 2.1, can easily be made available as oracle calls. The parties place their inputs in the appropriate registers  $\mathcal{A}_i^{\mathcal{O}, \text{in}}$  and  $\mathcal{B}_i^{\mathcal{O}, \text{in}}$ . The operation of the oracle  $\mathcal{O}_i$  is as follows: it measures the input registers  $\mathcal{A}_i^{\mathcal{O}, \text{in}}$  and  $\mathcal{B}_i^{\mathcal{O}, \text{in}}$  to force classical inputs. Then, it applies the appropriate classical ideal functionality on those classical inputs plus its classical internal state found in  $\mathcal{O}_{i-1}$ . This produces outputs for the parties and a new internal state. The outputs are placed in  $\mathcal{A}_i^{\mathcal{O}, \text{out}}$  and  $\mathcal{B}_i^{\mathcal{O}, \text{out}}$ . The new internal state of the oracle is placed in  $\mathcal{O}_i$ .

A two-party hybrid protocol for  $\mathcal{F} : L(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}) \rightarrow L(\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}})$  between parties  $\mathcal{A}$  and  $\mathcal{B}$  upon joint input state  $\rho_{\text{in}} \in D(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}} \otimes \mathcal{R})$  is defined as:

**Definition 2.5.** An  $m$ -turn two-party hybrid protocol  $\Pi_{\mathcal{F}}^{\mathcal{O}} = (\mathcal{A}, \mathcal{B}, \mathcal{O}, m)$  for  $\mathcal{F} : L(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}) \rightarrow L(\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}})$  is a  $m$ -turn two-party strategy with oracle calls, where  $\mathcal{A}_0 := \mathcal{A}_{\text{in}}$ ,  $\mathcal{B}_0 := \mathcal{B}_{\text{in}}$ ,  $\mathcal{A}_{m+1} := \mathcal{A}_{\text{out}}$ ,  $\mathcal{B}_{m+1} := \mathcal{B}_{\text{out}}$ , and where for all  $\rho_{\text{in}} \in D(\mathcal{A}_0 \otimes \mathcal{B}_0 \otimes \mathcal{R})$ ,  $\Delta([\mathcal{A} \otimes \mathcal{B}]^{\mathcal{O}}(\rho_{\text{in}}), (\mathcal{F} \otimes \mathbb{1}_{\mathcal{R}})(\rho_{\text{in}})) = 0$ . In the following, we often write simply two-party protocol to refer to a two-party hybrid protocol.

For  $i \in [0..m]$ , the joint state after turn  $i + 1$  in  $\Pi_{\mathcal{F}}^{\mathcal{O}}$  is denoted by  $[\mathcal{A} \otimes \mathcal{B}]_{i+1}^{\mathcal{O}}(\rho_{\text{in}}) := (\mathbb{1}_{L(\mathcal{B}_{i+1} \otimes \mathcal{A}_{i+1} \otimes \mathcal{R})} \otimes \mathcal{O}_{i+1})(\mathcal{A}_{i+1} \otimes \mathcal{B}_{i+1} \otimes \mathbb{1}_{L(\mathcal{R} \otimes \mathcal{O}_{i+1})})[\mathcal{A} \otimes \mathcal{B}]_i^{\mathcal{O}}(\rho_{\text{in}})$ , where  $[\mathcal{A} \otimes \mathcal{B}]_0^{\mathcal{O}}(\rho_{\text{in}}) := \rho_{\text{in}}$ , and  $[\mathcal{A} \otimes \mathcal{B}]_{m+1}^{\mathcal{O}}(\rho_{\text{in}}) := [\mathcal{A} \otimes \mathcal{B}]^{\mathcal{O}}(\rho_{\text{in}})$ .

### 3 Modeling Active Security

We start by extending the framework in [DNS10] to handle active security. Our model is very standard, defining security via simulation, but for completeness we describe and motivate the changes made in Appendix B.

Let  $\Pi_{\mathcal{F}}^{\mathcal{O}} = (\mathcal{A}, \mathcal{B}, \mathcal{O}, m)$  be a  $m$ -turn two-party hybrid protocol. Let  $\tilde{\mathcal{A}}$  and  $\tilde{\mathcal{B}}$  be adversaries in  $\Pi_{\mathcal{F}}^{\mathcal{O}}$ . We denote by  $[\tilde{\mathcal{A}} \otimes \mathcal{B}]^{\mathcal{O}}$  and  $[\mathcal{A} \otimes \tilde{\mathcal{B}}]^{\mathcal{O}}$  the resulting  $m$ -turn two-party strategies. We



will as usual define the security of such actively attacked protocols by comparing to a simulation. The simulation is basically an ideally secure evaluation of  $\mathcal{F}$ .

The ideally secure protocol for evaluating  $\mathcal{F}$  would be in a world where  $\mathcal{F}$  actually existed as an oracle—the parties would simply call this oracle. We can, however, not expect any protocol to be as secure as this, as for most protocols we cannot ensure that either both parties get the output or no party gets the output, known as fairness, which is ensured in the above ideal setting. We will therefore consider a setting where one of the parties learns its output first, and where the party learning it first, if corrupted, can prevent the other party from learning its output. We will only give the definition for the case where Alice learns first. Deriving the definition for the symmetric case where Bob learns first is trivial.

To avoid confusion between the parties in the real protocol and in the ideal protocol, we formulate the ideal protocol with parties Charleen,  $\mathcal{C}$ , and Dan,  $\mathcal{D}$ , taking the seats of Alice respectively Bob. So, we look at an ideal functionality  $\mathcal{F} : L(\mathcal{C}_{\text{in}} \otimes \mathcal{D}_{\text{in}}) \rightarrow L(\mathcal{C}_{\text{out}} \otimes \mathcal{D}_{\text{out}})$ , but where we keep a mental note reminding that  $\mathcal{C}_{\text{in}} := \mathcal{A}_{\text{in}}$ ,  $\mathcal{C}_{\text{out}} := \mathcal{A}_{\text{out}}$ ,  $\mathcal{D}_{\text{in}} := \mathcal{B}_{\text{in}}$  and  $\mathcal{D}_{\text{out}} := \mathcal{B}_{\text{out}}$ .

The ideal protocol for  $\mathcal{F}$  is then a 2-turn two-party hybrid protocol  $\Gamma_{\mathcal{F}} = (\mathcal{C}, \mathcal{D}, \mathcal{F}, 2)$ , which lets the parties query  $\mathcal{F}$  in turn 1, but only letting  $\mathcal{C}$  see her output in turn 1. In turn 2, Charleen then inputs a bit  $f$ , with  $f = 1$  indicating that Dan should receive his output and  $f = 0$  indicating that Dan should not receive his output. Dan will receive  $f$ , and if  $f = 1$  he will additionally be given his output. This is handled by the second oracle. Then the parties output whatever they received from the oracles. In the honest protocol, Charleen always inputs  $f = 1$ . We call this *the ideal protocol for evaluating  $\mathcal{F}$  without fairness for Dan*.

**The ideal protocol  $\Gamma_{\mathcal{F}} = (\mathcal{C}, \mathcal{D}, \mathcal{F}, 2)$  for  $\mathcal{F}$  without fairness for Dan:**

1. By convention we have  $\mathcal{C}_0 = \mathcal{C}_{\text{in}}$  and  $\mathcal{D}_0 = \mathcal{D}_{\text{in}}$  while  $\mathcal{O}_0$  is empty. In the ideal protocol, we let  $\mathcal{C}_1^{\mathcal{F}, \text{in}} \approx \mathcal{C}_0$  and  $\mathcal{D}_1^{\mathcal{F}, \text{in}} \approx \mathcal{D}_0$ .  $\mathcal{C}_1$  and  $\mathcal{D}_1$  simply send their respective input stored in  $\mathcal{C}_0$  and  $\mathcal{D}_0$  to  $\mathcal{C}_1^{\mathcal{F}, \text{in}}$  and  $\mathcal{D}_1^{\mathcal{F}, \text{in}}$  respectively. The oracle  $\mathcal{F}_1 : L(\mathcal{C}_1^{\mathcal{F}, \text{in}} \otimes \mathcal{O}_0 \otimes \mathcal{D}_1^{\mathcal{F}, \text{in}}) \mapsto L(\mathcal{C}_1^{\mathcal{F}, \text{out}} \otimes \mathcal{O}_1 \otimes \mathcal{D}_1^{\mathcal{F}, \text{out}})$ , we set  $\mathcal{C}_1^{\mathcal{F}, \text{out}} \approx \mathcal{C}_{\text{out}}$ ,  $\mathcal{O}_1 \approx \mathcal{D}_{\text{out}}$  and we let  $\mathcal{D}_1^{\mathcal{F}, \text{out}}$  be empty.  $\mathcal{F}_1$  simply applies the quantum operation  $\mathcal{F}$  to the inputs supplied by the parties and sends Charleen's output to Charleen in  $\mathcal{C}_1^{\mathcal{F}, \text{out}}$ . Dan's output is saved in the internal state  $\mathcal{O}_1$  of the oracle, giving Dan no output so far.
2. We let  $\mathcal{C}_2 \approx \mathcal{C}_1^{\mathcal{F}, \text{out}} \approx \mathcal{C}_{\text{out}}$  and we let  $\mathcal{C}_2^{\mathcal{F}, \text{in}}$  be a one qubit register, holding a qubit denoted  $|f\rangle$ . We let  $\mathcal{C}_2$  move the output from the oracle stored in  $\mathcal{C}_1^{\mathcal{F}, \text{out}}$  to  $\mathcal{C}_2$  and sets  $|f\rangle = |1\rangle$  in  $\mathcal{C}_2^{\mathcal{F}, \text{in}}$ .  $\mathcal{D}_2$  does nothing at this point.  
For the second oracle  $\mathcal{F}_2 : L(\mathcal{C}_2^{\mathcal{F}, \text{in}} \otimes \mathcal{O}_1 \otimes \mathcal{D}_2^{\mathcal{F}, \text{in}}) \mapsto L(\mathcal{C}_2^{\mathcal{F}, \text{out}} \otimes \mathcal{O}_2 \otimes \mathcal{D}_2^{\mathcal{F}, \text{out}})$ , we let  $\mathcal{C}_2^{\mathcal{F}, \text{out}}$  and  $\mathcal{O}_2$  be empty, and we let  $\mathcal{D}_2^{\mathcal{F}, \text{out}}$  be of the same dimension as  $\mathcal{D}_{\text{out}}$ , plus room for one qubit  $|a\rangle$ . It starts by measuring  $|f\rangle$  in the computational basis. If  $|f\rangle = |1\rangle$ , it then sets  $\mathcal{D}_2^{\mathcal{F}, \text{out}}$  to hold  $|a\rangle = |1\rangle$  along with the state in  $\mathcal{O}_1$ . If  $|f\rangle = |0\rangle$ , it sets  $\mathcal{D}_2^{\mathcal{F}, \text{out}}$  to hold  $|a\rangle = |0\rangle$  along with some fixed dummy state  $|\perp\rangle$  of the right dimension to fill  $\mathcal{C}_2^{\mathcal{F}, \text{out}}$ . When Charleen inputs  $f = 1$ , Dan will get his output. If  $f = 0$ , Dan gets no output, except a bit  $a = 0$  telling him that Charleen cheated him of his output (we say that Charleen *aborted* the computation).
3. We let  $\mathcal{C}_3 := \mathcal{C}_{\text{out}} \approx \mathcal{C}_2$  and  $\mathcal{D}_3 := \mathcal{D}_{\text{out}} \approx \mathcal{D}_2^{\mathcal{F}, \text{out}}$ .  $\mathcal{C}_3$  and  $\mathcal{D}_3$  simply output whatever they received from  $\mathcal{F}_1$  and  $\mathcal{F}_2$  in  $\mathcal{C}_2$  and  $\mathcal{D}_2^{\mathcal{F}, \text{out}}$  respectively, except qubit  $|a\rangle$  in the latter register. When the execution has aborted, Dan outputs a dummy state.

Consider the powers of a corrupted Charleen,  $\tilde{\mathcal{C}}$ . She might in the first round provide an alternative input for  $\mathcal{F}$ , possibly saving her original input, or a part thereof, in some ancilla. This is *input substitution*. Her choice of alternative input can only depend on her own original input, not that of Dan. This is *input independence*. Then she learns only the output of the oracle. This is *privacy*. After learning her own output she might then specify that Dan is not to learn his output.

This is the necessary *lack of fairness*. A corrupted Dan has similar powers, except that he cannot abort after seeing his output. We then say that a protocol is secure if it only allows attacks which could be mounted in the ideal protocol, i.e., it only allows the inevitable attack.

**Definition 3.1.** Let  $\Pi_{\mathcal{F}}^{\mathcal{O}} = (\mathcal{A}, \mathcal{B}, \mathcal{O}, m)$  be a two-party hybrid protocol for  $\mathcal{F} : \mathbb{L}(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}) \rightarrow \mathbb{L}(\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}})$ . Let  $\Gamma_{\mathcal{F}} = (\mathcal{C}, \mathcal{D}, \mathcal{F}, 2)$  be the ideal protocol for  $\mathcal{F}$  without fairness for Dan. Let  $\delta \in [0..1]$ . We say that  $\Pi_{\mathcal{F}}^{\mathcal{O}}$  is  $\delta$ -active secure without fairness for Bob, if for all adversaries  $\tilde{\mathcal{A}}$  and  $\tilde{\mathcal{B}}$  in  $\Pi_{\mathcal{F}}^{\mathcal{O}}$ , there exist adversaries  $\tilde{\mathcal{C}}$  and  $\tilde{\mathcal{D}}$  in  $\Gamma_{\mathcal{F}}$  with sizes polynomial in the sizes of  $\tilde{\mathcal{A}}$  and  $\tilde{\mathcal{B}}$  such that for all input states  $\rho_{\text{in}} \in \mathbb{D}(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}} \otimes \mathcal{R})$ ,

$$\Delta \left( [\tilde{\mathcal{A}} \otimes \tilde{\mathcal{B}}]^{\mathcal{O}}(\rho_{\text{in}}), [\tilde{\mathcal{C}} \otimes \tilde{\mathcal{D}}]^{\mathcal{F}}(\rho_{\text{in}}) \right) \leq \delta \text{ and } \Delta \left( [\mathcal{A} \otimes \tilde{\mathcal{B}}]^{\mathcal{O}}(\rho_{\text{in}}), [\mathcal{C} \otimes \tilde{\mathcal{D}}]^{\mathcal{F}}(\rho_{\text{in}}) \right) \leq \delta.$$

If a protocol is  $\delta$ -active secure for  $\delta = 0$ , then it is called perfectly active secure. If it is  $\delta$ -active secure for  $\delta$  negligible in some security parameter  $n$ , then it is called statistically active secure (in  $n$ ).

## 4 Securely Swaddling Wires to Ensure Ultimate Speciousness

The Clifford based QAS, described in Sect. 2.2, can be used to share the authentication of a quantum message in a way that allows any of the players, when helped by a TPC, to verify that a state has not been tampered with, and this without being able to get information on the encoded state. This primitive will be used extensively in our protocol to ensure that both players are *ultimately specious*. It relies on a string commitment scheme and secure TPCs provided as oracles, as described in Sect. 2.1.

The basic idea consists in swaddling each input wire  $w$  into a set of  $2n$  dummy wires where  $n$  belongs to Alice and  $n$  belongs to Bob. No party will be able to extract information about the state of the original wire from the swaddling. Moreover, any attempt to modify the state of the original wire will be detected except with negligible probability in  $n$ . The subprotocol  $\text{Swaddle}(w)$ , described below, uses two applications of the Clifford-based QAS (one on top of the other) in order to achieve this. Alice authenticates the state of her wire  $w$  and commits to her authentication key  $K_{w,a}$  using identifier  $id_A(w)$ . She then sends the resulting system to Bob who authenticates it using key  $K_{w,b}$  that he also commits upon using identifier  $id_B(w)$ . Bob sends back the resulting system to Alice allowing her to test the validity of the swaddling. The swaddling therefore uses a total of  $2n + 1$  wires, where Alice holds the *innermost* and Bob the *outermost* key of the resulting swaddling  $\mathfrak{s}(w)$  for wire  $w$ .

Consider now the test performed by Alice at Step 7. Let  $A$  and  $B$  be the Clifford operators corresponding to keys  $K_{w,a}$  and  $K_{w,b}$  respectively. Notice that Alice can easily test the authenticity of a swaddling when she holds the outermost authentication. She simply applies  $A^\dagger$  upon the  $2n + 1$  wires, measures her  $n$  dummy wires in the computational basis to verify that they are in state  $|0^n\rangle$ . She then re-applies  $A$  to the  $2n + 1$  wires of the swaddling. If Alice holds the innermost authentication of the swaddling, the testing procedure relies on the ideal functionality  $\text{id}_{\text{TEST}}$  defined as:

$$\text{id}_{\text{TEST}}((i_a, i_b), (i'_a, i'_b)) = \begin{cases} ((t, K'_{w,a}), K'_{w,b}) & \text{if } i_a = i'_a = id_A(w) \wedge i_b = i'_b = id_B(w) \text{ ,} \\ (\perp, \perp) & \text{otherwise ,} \end{cases}$$

where  $K'_{\mathbf{w},a}$ , and  $K'_{\mathbf{w},b}$  correspond to random Clifford operations  $A' \in \mathfrak{C}_{2n+1}$  and  $B' \in \mathfrak{C}_{n+1}$  respectively, and  $t$  corresponds to Clifford  $T \in \mathfrak{C}_{2n+1}$  subject to:

1.  $\text{id}_{\text{OPEN}}(\text{id}_A(\mathbf{w}), \text{id}_A(\mathbf{w})) = (\perp, K'_{\mathbf{w},a})$  and  $\text{id}_{\text{OPEN}}(\text{id}_B(\mathbf{w}), \text{id}_B(\mathbf{w})) = (K'_{\mathbf{w},b}, \perp)$ , and
2. If Alice holds the innermost key of  $\mathbf{w}$  then  $T = (\mathbb{1}_n \otimes B')(A^\dagger \otimes \mathbb{1}_n)B^\dagger$ .

**Swaddle( $\mathcal{W}$ ), with  $\mathcal{W} \subseteq \mathcal{A}$ :**

1. Alice initializes  $n \cdot \#\mathcal{W}$  dummy wires in state  $|0\rangle$ .
2. **For** each  $\mathbf{w} \in \mathcal{W}$ , Alice randomly chooses Clifford  $K_{\mathbf{w},a}$  on  $n + 1$  qubits, commits to it in  $\text{id}_{\text{SC}}(\text{id}_A(\mathbf{w}), K_{\mathbf{w},a}, \perp)$ , and applies it to  $\mathbf{w}$  and her dummies.
3. Alice **sends** these  $\#\mathcal{W} \cdot (n + 1)$  wires to Bob.
4. Bob initializes  $\#\mathcal{W} \cdot n$  dummy wires in the state  $|0\rangle$ .
5. **For** each  $\mathbf{w} \in \mathcal{W}$ , Bob randomly chooses a Clifford  $K_{\mathbf{w},b}$  on  $2n + 1$  qubits, commits to it in  $\text{id}_{\text{SC}}(\perp, (\text{id}_B(\mathbf{w}), K_{\mathbf{w},b}))$ , and applies it to the  $n + 1$  wires received from Alice as well as his own dummies.
6. Bob **sends** all  $\#\mathcal{W} \cdot (2n + 1)$  wires back to Alice. Let  $\mathfrak{s}(\mathbf{w})$  denote the resulting swaddling of  $\mathbf{w} \in \mathcal{W}$ .
7. **For** each  $\mathbf{w} \in \mathcal{W}$ , Alice **calls**  $\text{TestSwaddling}(\mathfrak{s}(\mathbf{w}))$ .

Condition 1 ensures that the authentication keys with identifiers  $\text{id}_A(\mathbf{w})$  and  $\text{id}_B(\mathbf{w})$  have been updated to hold values  $K'_{\mathbf{w},a}$  and  $K'_{\mathbf{w},b}$  respectively. Condition 2 makes sure that the state  $|\sigma\rangle$  of a valid swaddling  $\mathfrak{s}(\mathbf{w})$  satisfies  $T|\sigma\rangle = |0^n\rangle \otimes B'(|\varphi\rangle \otimes |0^n\rangle)$ , where  $|\varphi\rangle$  is the logical state of  $\mathfrak{s}(\mathbf{w})$ . Notice that Alice gets no information about Bob's Clifford  $B'$  if she had no information about  $B$  to start with. The security of  $\text{Swaddle}(\mathbf{w})$  will be established in Appendix G.

**TestSwaddling( $\mathfrak{s}(\mathbf{w})$ ), with Alice doing the testing:**

1. **If** Alice holds the outermost authentication key  $K_{\mathbf{w},a}$  **then**:
  - Alice applies  $A^\dagger$  on  $\mathfrak{s}(\mathbf{w})$ , where  $A$  is the Clifford corresponding to string  $K_{\mathbf{w},a}$ ,
  - Alice tests that her dummies are together in state  $|0^n\rangle$ . **Otherwise**, she **aborts**.
  - Alice re-applies  $A$  on the  $2n + 1$  qubits of the swaddling.
2. **Else** they **call**  $((t, K'_{\mathbf{w},a}), K'_{\mathbf{w},b}) = \text{id}_{\text{TEST}}(\text{id}_A(\mathbf{w}), \text{id}_B(\mathbf{w}), (\text{id}_A(\mathbf{w}), \text{id}_B(\mathbf{w})))$ ,
  - **If** one party gets  $\perp$  in  $\text{id}_{\text{TEST}}$  **then abort**.
  - Alice applies  $T$  on  $\mathfrak{s}(\mathbf{w})$  where  $T$  is the Clifford operator corresponding to string  $t$ .
  - Alice tests that her dummies are together in state  $|0^n\rangle$ . **Otherwise**, she **aborts**.
  - Alice applies  $A'$  to the  $2n + 1$  qubits of the swaddling (note that  $A'$  and  $B'$  are committed upon with identifiers  $\text{id}_A(\mathbf{w})$  and  $\text{id}_B(\mathbf{w})$  respectively).

At the end of  $\text{TestSwaddling}(\mathfrak{s}(\mathbf{w}))$ , Alice holds the outermost authentication key  $K'_{\mathbf{w},a}$  of the resulting swaddling. Bob can do the testing by the exact same procedure provided the roles of Alice and Bob are reversed in  $\text{TestSwaddling}(\mathfrak{s}(\mathbf{w}))$ . The security of these procedures will be discussed in Sect. 5.5. Intuitively, we expect that  $\text{TestSwaddling}(\mathfrak{s}(\mathbf{w}))$  allows parties to test that a swaddling  $\mathfrak{s}(\mathbf{w})$  has not been tampered with. It will also be used to test that each party transforms a swaddling  $\mathfrak{s}(\mathbf{w})$  the way they should during the execution of the protocol. Notice that no information about the logical state of wire  $\mathbf{w}$  leaks to any party in  $\text{Swaddle}(\mathbf{w})$  and  $\text{TestSwaddling}(\mathfrak{s}(\mathbf{w}))$  since statistically secure authentication must also encrypt  $\mathbf{w}$  [BCG<sup>+</sup>02].

At the end of our protocol, each party will be asked to verify that the other party's registers upon which the circuit is evaluated are all in the states they should be. This subprotocol is called  $\text{TestAllSwaddlings}$ . It simply consists in the execution of  $\text{TestSwaddling}(\mathfrak{s}(\mathbf{w}))$  for all wires  $\mathbf{w}$  held by each party.

**TestAllSwaddlings:**

1. For all wires  $w$  held by Bob, Bob sends  $s(w)$  to Alice before Alice and Bob run  $\text{TestSwaddling}(s(w))$  with Alice doing the testing.
2. All swaddles  $s(w)$  received by Alice are sent back to Bob before Alice and Bob run  $\text{TestSwaddling}(s(w))$  with Bob doing the testing.
3. For all wires  $w$  held by Alice, Alice sends  $s(w)$  to Bob before Alice and Bob run  $\text{TestSwaddling}(s(w))$  with Bob doing the testing.
4. All swaddles  $s(w)$  received by Bob are sent back to Alice before Alice and Bob run  $\text{TestSwaddling}(s(w))$  with Alice doing the testing.

The following subprotocol implements the obvious way the openings of all committed Clifford operators, thereby allowing Alice and Bob to get the final state of the computation. Notice that since Alice learns Bob's secret keys before she unveils her own, our protocol will lack fairness for Bob.

**OpenAllSwaddlings:**

1. **For** each wire  $w \in \mathcal{A}_{\text{out}}$ , Bob reveals to Alice his committed secret key encrypting  $w$ , with identifier  $id_B(w)$ , by **calling**  $\text{id}_{\text{OPEN}}(id_B(w), id_B(w))$ . **If** Alice or Bob get  $\perp$  from the ideal functionality **then abort**.
2. **For** each wire  $w \in \mathcal{B}_{\text{out}}$ , Alice reveals to Bob her committed secret key encrypting  $w$ , with identifier  $id_A(w)$ , by **calling**  $\text{id}_{\text{OPEN}}(id_A(w), id_A(w))$ . **If** Alice or Bob get  $\perp$  from the ideal functionality **then abort**.

## 5 Description of the protocol

We first start by defining the various spaces on which (the honest) Alice and Bob will be working. The circuit that they want to execute acts on some wires in Alice's possession and some in Bob's possession, and these will be swaddled as described above. We will denote by  $\mathcal{A}_u, \mathcal{B}_u$  the spaces corresponding to Alice and Bob's unswaddled wires, reserving  $\mathcal{A}$  and  $\mathcal{B}$  for the actual swaddled wires.

The core idea of our protocol is the same as in [DNS10]: we first initialize all the qubits by swaddling them, we then perform each gate from the circuit one after the other on the authenticated data. Hence, we need to give subprotocols for the initialization as well as for each of the gates in our universal set.

For all Clifford gates (i.e., all gates in  $\mathcal{UG}$  except for R), the subprotocols are fairly simple: we use classical two-party computation to reveal a Clifford operation that executes the gate while updating the encryption key; the revealed Clifford then looks uniformly distributed and independent of everything else.

Implementing the R-gate is more involved. We use ideas from fault-tolerant computation, where this gate is implemented by doing gate teleportation via a so-called *magic state*: one prepares a special state (namely  $|M\rangle = \frac{1}{\sqrt{2}}(|0\rangle + e^{i\pi/4}|1\rangle)$ ) and then use a teleportation-like circuit, which itself requires only Clifford gates and measurements, to execute the gate. The problem is then reduced to that of producing this magic state, which can be done by a distillation process. The distillation process that we use is exactly the one considered in [BK05] (where  $|M\rangle$  is an “ $H$ -type magic state” in their language); a description of it can also be found in Appendix D.

### 5.1 Main Protocol

We now give the full description of our protocol, denoted  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{O}'}$  =  $(\mathcal{A}, \mathcal{B}, \mathcal{O}', m)$ , allowing to evaluate the CPTP map  $\mathcal{F} : L(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}) \mapsto L(\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}})$  upon joint input state  $\rho_{\text{in}} \in$

$D(\mathcal{R} \otimes \mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}})$ . The oracle list needed to run the protocol in the hybrid model is provided implicitly in Sect. 5. The operations performed at each step by  $\mathcal{A}$  and  $\mathcal{B}$  are described informally as the instructions of Alice and Bob respectively. We will view  $\mathcal{F}$  as being implemented by a quantum circuit acting on  $\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}$  together with ancillas  $\mathcal{A}_a$  and  $\mathcal{B}_a$  initialized in state  $|0\rangle$  (we shall explain below how to test that ancillas are really in state  $|0\rangle$ ). At the end of the circuit, some of these wires become part of the outputs  $\mathcal{A}_{\text{out}}$  and  $\mathcal{B}_{\text{out}}$ , and the rest are part of the environment that will remain encrypted ( $\mathcal{A}_e$  and  $\mathcal{B}_e$ ). In the following, we denote by  $U_{\mathcal{F}} \in U(\mathcal{A}_{\text{in}} \otimes \mathcal{A}_a \otimes \mathcal{B}_{\text{in}} \otimes \mathcal{B}_a)$  the unitary transform implemented by the circuit that satisfies, for all mixed quantum state  $\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}}$ ,

$$\mathcal{F}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}}) = \text{tr}_{\mathcal{A}_e\mathcal{B}_e} \left( U_{\mathcal{F}} \rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}} \otimes |0\rangle\langle 0|^{\mathcal{A}_a} \otimes |0\rangle\langle 0|^{\mathcal{B}_a} U_{\mathcal{F}}^\dagger \right),$$

where  $\mathcal{A}_e \subseteq \mathcal{A}_{\text{in}} \otimes \mathcal{A}_a$ ,  $\mathcal{B}_e \subseteq \mathcal{B}_{\text{in}} \otimes \mathcal{B}_a$ .

Let  $G_1, G_2, \dots, G_{\ell(n)}$  be an enumeration of all gates of the circuit for  $\mathcal{F}$  where  $\ell(n)$  is polynomial in  $n$ , and  $G_i$  is executed before  $G_{i+1}$ . This protocol calls a number of subprotocols described next and in Appendix E.

**Protocol  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{G}'}$  for the evaluation of  $\mathcal{F}$  upon joint input  $\rho_{\text{in}} \in D(\mathcal{R} \otimes \mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}})$ :**

1. Alice and Bob **run Initialization**,
2. **For**  $i = 1 \dots \ell(n)$ :
  - **If**  $G_i$  is a one-bit Clifford gate applied to wire  $w$  **then** Alice and Bob **call** `OneQubitClifford( $G_i, s(w)$ )`,
  - **If**  $G_i$  is a CNOT-gate applied to control wire  $w_c$  and  $w_t$  **then** Alice and Bob **call** `CNOT( $s(w_c), s(w_t)$ )`,
  - **If**  $G_i$  is an R-gate applied to wire  $w$  **then** Alice and Bob **call** `RGate( $s(w)$ )`,
3. Alice and Bob **call** `TestAllSwaddlings`,
4. Alice and Bob **call** `OpenAllSwaddlings`,
5. Alice and Bob decrypt the swaddlings for all wires in  $\mathcal{A}_{\text{out}}$  and  $\mathcal{B}_{\text{out}}$  using the keys received by the other party together with their own.

## 5.2 Subprotocols

The `Initialization` subprotocol prepares all swaddlings required during the evaluation of the circuit. In addition to swaddling all wires holding the qubits upon which the circuit acts, some ancillary states also have to be swaddled: one wire in the magic state and one wire in state  $|0\rangle$  per R-gate in the circuit.

**Initialization:**

1. **For** each R-gate being applied to a wire in  $\mathcal{A}$ , Alice adds one additional wire  $m_{a,i}$  initialized in state  $|M\rangle = \frac{1}{\sqrt{2}}(|0\rangle + e^{i\pi/4}|1\rangle)$  to  $\mathcal{A}$ .
2. **For** each R-gate being applied to a wire in  $\mathcal{B}$ , Bob adds an additional wire  $m_{b,i}$  initialized in state  $|M\rangle$  to  $\mathcal{B}$ .
3. **For** each R-gate being applied to a wire in  $\mathcal{B}$ , Alice adds an additional wire  $c_{a,i}$  initialized in state  $|0\rangle$  to  $\mathcal{A}$ .
4. **For** each R-gate being applied to a wire in  $\mathcal{A}$ , Bob adds an additional wire  $c_{b,i}$  initialized in state  $|0\rangle$  to  $\mathcal{B}$ .
5. Alice **calls** `Swaddle( $\mathcal{A}$ )`.
6. Bob **calls** `Swaddle( $\mathcal{B}$ )`.
7. **For** every wire  $a \in \mathcal{A}_a$  or added in step 3 above, Alice **calls** `VerifyAncilla( $s(a)$ )`.
8. **For** every wire  $a \in \mathcal{B}_a$  or added in step 4 above, Bob **calls** `VerifyAncilla( $s(a)$ )`.
9. **For** each magic wire  $m \in \mathcal{A}$ , Alice **calls** `DistillMagic( $s(m)$ )`.
10. **For** each magic wire  $m \in \mathcal{B}$ , Bob **calls** `DistillMagic( $s(m)$ )`.

Subprotocols `VerifyAncilla( $s(a)$ )` and `DistillMagic( $s(m)$ )` are described in Appendices E and D. The following subprotocols evaluates any one-qubit Clifford, the CNOT-gate, respectively the R. The idea behind the one-qubit protocol is to use TPC to compute the gate *through* the

authentication. The idea in the CNOT-subprotocol is to take the two swaddled qubits involved and turn them into one big swaddling. Then, performing a CNOT on them is no different from performing any other Clifford gate on the whole block. We then separate the big swaddling back into its components.

To express the TPCs needed for these subprotocols, we will need to define some more ideal functionalities. For the one-qubit Clifford gates, we will define  $\text{id}_{1\text{QC}}$ , which simply updates the commitment to the inner key of the swaddling: a call to  $\text{id}_{1\text{QC}}((i, G), (i', G'))$ , if  $i = i' = \text{id}_A(\mathbf{w})$  for some wire  $\mathbf{w}$  for which Alice has the inner swaddling and  $G = G'$  is some one-qubit Clifford, results in the commitments to the key  $K_{\mathbf{w},a}$  to be updated to commitments to  $K_{\mathbf{w},a}G^\dagger$  (and likewise when Bob holds the inner key).

The subprotocols are then as follows:

**OneQubitClifford( $C, \mathfrak{s}(\mathbf{w})$ ):**

1. If Alice has the inner swaddling of  $\mathbf{w}$ , Alice and Bob call  $\text{id}_{1\text{QC}}((\text{id}_A(\mathbf{w}), C), (\text{id}_A(\mathbf{w}), C))$ , otherwise they call  $\text{id}_{1\text{QC}}((\text{id}_B(\mathbf{w}), C), (\text{id}_B(\mathbf{w}), C))$ .

**CNOT( $\mathfrak{s}(\mathbf{w}_c), \mathfrak{s}(\mathbf{w}_t)$ ) with wire  $\mathbf{w}_c \in \mathcal{A}$  as control and  $\mathbf{w}_t$  as target:**

1. Bob **sends**  $\mathfrak{s}(\mathbf{w}_t)$  to Alice if Bob holds it, in which case Alice **calls**  $\text{TestDummies}(\mathfrak{s}(\mathbf{w}_t))$ .
2. Alice performs a randomly selected Clifford  $C$  on  $4n + 2$  qubits jointly on  $\mathfrak{s}(\mathbf{w}_c)$  and  $\mathfrak{s}(\mathbf{w}_t)$ .
3. Alice **sends**  $\mathfrak{s}(\mathbf{w}_c)$  and  $\mathfrak{s}(\mathbf{w}_t)$  to Bob; Bob **calls**  $\text{TestDummies}$  jointly on them.
4. Alice and Bob perform a TPC whose outcome tells Bob to perform a Clifford unitary  $C' = (K'_c \otimes K'_t)(\text{CNOT})K^\dagger$  where  $K$  is the key of the swaddling at this point, and  $K'_c$  and  $K'_t$  are randomly-chosen Cliffords that become the new key.
5. Bob **sends**  $\mathfrak{s}(\mathbf{w}_c)$  (and  $\mathfrak{s}(\mathbf{w}_t)$  if she held this one too) to Alice. Alice **calls**  $\text{TestDummies}$  on them.

**RGate( $\mathfrak{s}(\mathbf{w})$ ) with  $\mathbf{w} \in \mathcal{A}$ :**

- 1.
2. Alice performs a swaddled CNOT with  $\mathbf{m}_w$  as a control, and  $\mathbf{w}$  as target.
3. Alice **sends**  $\mathfrak{s}(\mathbf{w})$  to Bob. Bob **calls**  $\text{TestSwaddling}(\mathfrak{s}(\mathbf{w}))$  on it.
4. Alice **calls**  $\text{Measure}(\mathfrak{s}(\mathbf{w}))$ .
5. Alice and Bob perform a TPC whose result is a Clifford which, if both measurement results were zero, updates the key, and if both measurement results were one, performs a swaddled  $e^{i\pi/4}XP^\dagger$  and then updates the key. If the measurement results differ, then they abort.
6. Alice relabels  $\mathbf{m}_{\alpha,i}$  to  $\mathbf{w}$ .

The RGate subprotocol performs the R-gate using gate teleportation[GC99b,GC99a] via the magic state similarly to the fault-tolerant version of the R-gate introduced in [Got10,Sho96]. To perform an R-gate on a wire  $\mathbf{w}$  via gate teleportation, we would first perform a CNOT from the magic state to  $\mathbf{w}$ , and then measure  $\mathbf{w}$  in the computational basis. If the answer is 0, we do nothing, and if it is 1, we perform  $e^{i\pi/4}XP^\dagger$  on the former magic state, that we then rename  $\mathbf{w}$ .  $\text{Measure}(\mathfrak{s}(\mathbf{w}))$  is described in Appendix E.

### 5.3 Security of the subprotocols

We will now show that the protocol described in the previous section has the following property: any adversary which deviates significantly from the protocol will be caught cheating with high probability. The general strategy will be as follows. For the initialization, any adversary will be forced to input something into the protocol, and will end up with some state that is properly

swaddled. Then, for every other protocol step, we will assume that at the beginning, the inputs are properly swaddled, and will aim to show that after the protocol step is done, we are once again left with a correct swaddling of the data, to which the correct operation has been applied. Furthermore, at every step, any deviation from the protocol will be essentially equivalent to an attack on the authentication scheme, which means that an adversary's chances of succeeding in changing the state without getting caught will be negligible in the number of dummy wires per qubit.

#### 5.4 Some additional definitions

Before we start, we will find it convenient to introduce some additional notation. From now on,  $\rho_0$  will consist of  $\rho_{\text{in}}$  augmented with all additional qubits introduced in the Initialization phase: the additional ancillas, and the dummy wires Alice and Bob use for swaddling. Furthermore, we will denote by  $\mathcal{K}_a$  and  $\mathcal{K}_b$  systems which represent Alice's and Bob's current key, respectively. These should be thought of as being part of the inner state of the TPC ideal functionality, and therefore cannot be changed at will.  $\mathcal{A}$  represents all other systems at Alice,  $\mathcal{B}$  represents all systems in Bob's possession, and  $\mathcal{R}$  is a system that includes everything else and ensures that the total state is pure.

In the sequel, we will call a *step* an execution of any of the subprotocols listed above; each of these steps consists of multiple *turns* in the sense of Section 2.3: the state at turn  $i$  consists of the state before the  $i$ th use of an oracle in the protocol (either a communication oracle or a classical computation oracle).

We will denote by  $C_{a,s}$  the operation that encodes Alice's qubits according to the keys stored in the TPC ideal functionalities at turn  $s$  in the protocol, and likewise for Bob's encoding operation  $C_{b,s}$ . See Appendix F for more precise definitions. Furthermore, we will denote by  $[\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \overline{\text{E-ABORT}}]_s^{\mathcal{O}'}$  ( $\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}}$ ) the global state of the protocol at turn  $s$  conditioned on the fact that the protocol doesn't abort before `TestAllSwaddlings` is completed. Note that this state is not normalized; its trace corresponds to the probability of not aborting before the end.

**Definition 5.1 (Forcing).** *We will say that the protocol is  $\mathcal{G}^{\mathcal{A}\mathcal{B}}$ -forcing for  $\tilde{\mathcal{A}}$  at turn  $s$  with initial operation  $\tilde{\mathcal{E}}_0^{\mathcal{A}}$  if there exists a final operation  $\tilde{\mathcal{E}}^{\mathcal{A}}$  such that:*

$$\Delta \left( [\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \overline{\text{E-ABORT}}]_s^{\mathcal{O}'}, \tilde{\mathcal{E}} \circ C_{b,s} \circ \mathcal{G} \circ \tilde{\mathcal{E}}_0 \left( \rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}} \right) \right) \leq \text{negl}(n) ,$$

where  $\tilde{\mathcal{E}}_0$  is a completely positive, trace non-increasing map that acts only on qubits in Alice's possession at the beginning of the protocol: her own input qubits, the dummies she inputs in the swaddling, and the various ancillas that she adds.

In other words, if the protocol is  $\mathcal{G}$ -forcing for Alice at some turn  $s$ , then, if the protocol doesn't abort early, regardless of what she tries to do, it will be essentially equivalent to changing her input (using the initial operation  $\tilde{\mathcal{E}}_0$ ), executing the operation  $\mathcal{G}$  (which will turn out to be the circuit that is supposed to be executed up to turn  $s$ ), swaddling the result, and then doing an arbitrary operation on her share alone, represented by  $\tilde{\mathcal{E}}$ .

**Definition 5.2.** *We say that a subprotocol is  $\mathcal{G}$ -forcing if the protocol is  $\mathcal{G} \circ \mathcal{G}_0$ -forcing at the end of the subprotocol given that it was  $\mathcal{G}_0$ -forcing at the beginning.*

**Definition 5.3 (Hiding).** We will say that the protocol is hiding for  $\tilde{\mathcal{A}}$  at turn  $s$ , if, for all input states  $\rho_{\text{in}}^{A_{\text{in}}B_{\text{in}}\mathcal{R}_{\text{in}}}$ , we have that

$$\Delta \left( \text{tr}_{\mathcal{B}} \left( [\tilde{\mathcal{A}} \otimes \mathcal{B}]_s^{\mathcal{O}'} (\rho_{\text{in}}^{A_{\text{in}}B_{\text{in}}\mathcal{R}_{\text{in}}}) \right), \text{tr}_{\mathcal{B}} \left( [\tilde{\mathcal{A}} \otimes \mathcal{B}]_s^{\mathcal{O}'} (\rho_{\text{in}}^{A_{\text{in}}\mathcal{R}_{\text{in}}} \otimes |0\rangle\langle 0|^{B_{\text{in}}}) \right) \right) \leq \text{negl}(n) .$$

Hence, the protocol is hiding if the state seen by a dishonest Alice is independent of Bob’s input. Note that this in particular implies that  $\tilde{\mathcal{A}}$ ’s action at turn  $s$  is necessarily independent of Bob’s input.

## 5.5 Proving hidingness and forcingness

To construct the simulator needed to prove security, we will need to show two things. First, we will have to show that before the keys are revealed, the cheater’s internal state can be produced by running the protocol internally with a dummy input from the honest party; this will follow from the fact that the Clifford QAS is a secure encryption scheme. Second, we will have to show that after the keys are revealed, the correct circuit has been applied. These two properties correspond to the “hiding” and “forcing” properties defined in the previous section, and proving these for our protocol will be the focus of this section. The fact that the protocol is hiding simply follows from the security of the Clifford QAS as an encryption scheme; we state this as a lemma below. To prove forcingness, the usual trick will be to assume that we pass every call to `TestDummies` in the protocol (since we only need to look at the no-early-abort case) and use the security definition of the Clifford QAS (Definition 2.2) to show that the dishonest party’s attack can be represented by a completely positive, trace non-increasing map  $\mathcal{U}^{\text{acc}}$  that acts only on his/her other systems (i.e. the ones that were not involved in the `TestDummies`). If the adversary decides to try to break the QAS at this step, this  $\mathcal{U}^{\text{acc}}$  will simply decrease the trace to reflect the probability of abortion. We prove the forcingness of the various subprotocols in Appendix G, and simply summarize the end result as Lemma 5.5 below.

**Lemma 5.4.** For every turn before `OpenAllSwaddlings`, the protocol is hiding for every adversary  $\tilde{\mathcal{A}}$  and  $\tilde{\mathcal{B}}$ .

*Proof.* During this phase of the protocol, all  $\tilde{\mathcal{A}}$  ever gets from Bob (and  $\tilde{\mathcal{B}}$  from Alice) is encrypted in a QAS, whose key is managed by the classical two-party computations. Hence, this follows directly from the security of the Clifford QAS (see Section 2.2 and Appendix C).

**Lemma 5.5.** All subprotocols are  $\mathcal{G}$ -forcing for any  $\tilde{\mathcal{A}}$  and  $\tilde{\mathcal{B}}$ , where  $\mathcal{G}$  is the operation performed by the subprotocol.

## 6 Proving Active Security

In this section, we prove active security of  $\hat{\Pi}_{\mathcal{F}}^{\mathcal{O}'}$  by providing a simulator  $\tilde{\mathcal{C}}$ . We prove that for any shared input state, when  $\tilde{\mathcal{C}}$  interacts with  $\mathcal{D}$ , it reproduces the state generated by  $\tilde{\mathcal{A}}$  interacting with  $\mathcal{B}$  on that same input state.

At the beginning,  $\tilde{\mathcal{C}}$  simulates the view of any execution between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}$  that aborts before all swaddlings are opened. This part of the simulation is described and analyzed in the next subsection. In Sect. 6.2, we provide and analyze the simulation when no early aborting occurs between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}$ .



## 6.1 Simulation of Early Aborting Executions

Here is the rough description of  $\tilde{\mathcal{C}}$ 's operations for the start of the simulation. These steps allow to simulate any execution aborting in the real world before the adversary  $\tilde{\mathcal{A}}$  receives back all her swaddlings from  $\mathcal{B}$  in `TestAllSwaddlings`. Let  $s^*$  be the turn in  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{O}'}$  at which this transmission is received by  $\tilde{\mathcal{A}}$ . The simulator  $\tilde{\mathcal{C}}$  runs  $\tilde{\mathcal{A}}$  as a subroutine using an internal copy of  $\mathcal{B}$ 's instructions run upon a (or any) dummy input state  $|0\rangle^{\mathcal{B}_{\text{in}}}$ . We denote by  $\mathcal{B}^*$  the simulated  $\mathcal{B}$  on a dummy input run internally in  $\tilde{\mathcal{C}}$ . In order to distinguish the registers of the real  $\mathcal{B}$  from the ones of  $\mathcal{B}^*$ , register  $\mathcal{B}$  of  $\mathcal{B}$  will be denoted by  $\mathcal{B}^*$  when used by  $\mathcal{B}^*$ . Consider the simulator  $\tilde{\mathcal{C}}$  starting its simulation as follows:

1.  $\tilde{\mathcal{C}}$  simulates the protocol execution between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$  holding the dummy input state  $|0^n\rangle^{\mathcal{B}_{\text{in}}^*}$  until turn  $s^*$  is reached. Remember that  $\tilde{\mathcal{C}}$  simulates all oracle calls during the execution between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$ . In particular  $\tilde{\mathcal{C}}$  stores all queries to classical oracles made by  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$ .
2. If the execution **aborts** at any point before  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$  reach turn  $s^*$  then  $\tilde{\mathcal{C}}$  **aborts** as well after having deleted all its internal registers except the ones held by  $\tilde{\mathcal{A}}$ .

At Step 2,  $\tilde{\mathcal{C}}$  and  $\mathcal{B}^*$  **aborts** with the same probability as a real world execution aborts prior to turn  $s^*$ . The simulation reproduces  $\tilde{\mathcal{A}}$ 's view in the real world when such an aborting execution occurs. This follows from the fact that up to that point,  $\mathcal{B}^*$  never provides any information about its input state. We now prove that the simulation of an early abort is correct. We only provide the proof that the simulator does the job when the adversary  $\tilde{\mathcal{A}}$  is impersonating  $\mathcal{A}$ . The proof for an adversary  $\tilde{\mathcal{B}}$  impersonating  $\mathcal{B}$  follows the same lines.

We define  $\text{E-ABORT}_s$  as the event consisting in an execution between two parties aborting at turn  $s < s^*$ . Let  $[\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}' }(\rho_{\text{in}})$  be denoting the the joint state of an execution that aborts at turn  $s$  whenever  $\tilde{\mathcal{A}}$  and  $\mathcal{B}$  are interacting upon joint input state  $\rho_{\text{in}}$ . Let  $[\tilde{\mathcal{C}} \otimes \mathcal{D} \wedge \text{E-ABORT}_s]^{\mathcal{F}}(\rho_{\text{in}})$  be denoting an execution between  $\tilde{\mathcal{C}}$  and  $\mathcal{D}$  (in the ideal world) upon joint input state  $\rho_{\text{in}}$  aborting while  $\tilde{\mathcal{C}}$  is simulating  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$  at turn  $s$  upon input state  $\text{tr}_{\mathcal{B}_{\text{in}}}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}}) \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}}$ . States  $[\tilde{\mathcal{C}} \otimes \mathcal{D} \wedge \text{E-ABORT}_s]^{\mathcal{F}}(\rho_{\text{in}})$  and  $[\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}' }(\rho_{\text{in}})$  are not normalized,  $\text{tr}([\tilde{\mathcal{C}} \otimes \mathcal{D} \wedge \text{E-ABORT}_s]^{\mathcal{F}}(\rho_{\text{in}}))$  and  $\text{tr}([\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}' }(\rho_{\text{in}}))$  are the probabilities that  $[\tilde{\mathcal{C}} \otimes \mathcal{D}]^{\mathcal{F}}(\rho_{\text{in}})$  and  $[\tilde{\mathcal{A}} \otimes \mathcal{B}]^{\mathcal{O}' }(\rho_{\text{in}})$  *aborts* at step  $s$  respectively. We have,

**Lemma 6.1 (Early abort).** *Let  $\tilde{\mathcal{A}}$  be an adversary in hybrid protocol  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{O}' } = (\mathcal{A}, \mathcal{B}, \mathcal{O}', m)$  for  $\mathcal{F} : \mathcal{L}(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}) \rightarrow \mathcal{L}(\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}})$ . Let  $s^*$  be the turn in  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{O}' }$  at which  $\mathcal{B}$  returns all of  $\tilde{\mathcal{A}}$ 's swaddlings in `TestAllSwaddlings` and let  $0 \leq s < s^*$ . Then, there exists an adversary  $\tilde{\mathcal{C}}$  in  $\Gamma_{\mathcal{F}}$  (polysize in the size of  $\tilde{\mathcal{A}}$ ) such that for any  $\rho_{\text{in}} \in \mathcal{D}(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}} \otimes \mathcal{R})$ ,*

$$\Delta \left( [\tilde{\mathcal{C}} \otimes \mathcal{D} \wedge \text{E-ABORT}_s]^{\mathcal{F}}(\rho_{\text{in}}), [\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}' }(\rho_{\text{in}}) \right) \leq \text{negl}(n) .$$

*By symmetry, the same is also true with respect to adversaries  $\tilde{\mathcal{B}}$  in  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{O}' }$  and  $\tilde{\mathcal{D}}$  in  $\Gamma_{\mathcal{F}}$ .*

*Proof (sketch).* The proof proceeds along the same lines than in [DNS10]. Since all communications between parties are statistically encrypted, an execution between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}$  upon joint input state  $\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}}$  is identical to an execution with joint input state  $\rho_{\text{in}}$  from the adversary's

point of view. Remember that when  $\tilde{\mathcal{C}} \otimes \mathcal{D}$  and  $\tilde{\mathcal{A}} \otimes \mathcal{B}$  abort at turn  $s < s^*$ ,  $\mathcal{D}$  and  $\mathcal{B}$  output a special state  $|\perp\rangle\langle\perp|$ . We have,

$$\begin{aligned} & \Delta \left( [\tilde{\mathcal{C}} \otimes \mathcal{D} \wedge \text{E-ABORT}_s]^{\mathcal{F}}(\rho_{\text{in}}), [\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}'_1}(\rho_{\text{in}}) \right) \\ &= \Delta \left( [\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}'_1}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}}), [\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}'_1}(\rho_{\text{in}}) \right) \end{aligned} \quad (2)$$

$$\begin{aligned} &= \Delta \left( \text{tr}_{\mathcal{B}} \left( [\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}'_1}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}}) \right), \right. \\ & \quad \left. \text{tr}_{\mathcal{B}} \left( [\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \text{E-ABORT}_s]^{\mathcal{O}'_1}(\rho_{\text{in}}) \right) \right) \end{aligned} \quad (3)$$

$$\leq \text{negl}(n) , \quad (4)$$

where (2) is a consequence of how the simulator is constructed, and (3) follows from the fact that in such executions, the output of the honest player is  $|\perp\rangle\langle\perp|$ . Equation (4) follows from the hiding property at every turn stated in Lemma 5.4. Notice that some measurements can be applied to some swaddled ancillas (initially in state  $|0\rangle$ ) or magic states while performing `VerifyMagic` or evaluating an `R`-gate. None of these measurements can reveal anything non-negligible about  $\mathcal{B}$ 's input state. For the evaluation of an `R`-gate, the outcome of measurements are independent of  $\mathcal{B}$ 's input state. For the magic state distillation process, measurements are only applied to candidate magic states also independent of  $\mathcal{B}$ 's input state.

To prove the statement for adversary  $\tilde{\mathcal{B}}$  interacting with  $\mathcal{A}$  in  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{O}'_1}$ , we define  $s^*$  to be the turn at which  $\tilde{\mathcal{B}}$  receives all his swaddlings from  $\mathcal{A}$ . The simulator  $\tilde{\mathcal{D}}$  for executions aborting early proceeds the same way than  $\tilde{\mathcal{C}}$  except that  $\tilde{\mathcal{B}}$  replaces  $\tilde{\mathcal{A}}$  while  $\mathcal{A}^*$  replaces  $\mathcal{B}^*$ .  $\square$

## 6.2 Simulation of Executions that do not Abort Early

It remains to simulate the execution from turn  $s^*$  until the end. In the simulated world, if  $\tilde{\mathcal{C}}$  reaches the end of Step 2 then the output state  $\mathcal{F}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|)$  can be recovered. The reason being that at turn  $s^*$  (taking place before `OpenAllSwaddlings`), all swaddlings have been tested by  $\mathcal{B}^*$  in `TestAllSwaddlings`.  $\tilde{\mathcal{C}}$  can get the output state since it knows all keys allowing to decrypt all logical wires, and all these wires have not been tampered with. The simulation then works along the same lines than in [DNS10]. At turn  $s^*$ ,  $\tilde{\mathcal{C}}$  is simulating  $\mathcal{B}^*$ 's quantum transmission of all swaddlings belonging to  $\tilde{\mathcal{A}}$  back to  $\tilde{\mathcal{A}}$ . These swaddlings have been successfully tested by  $\mathcal{B}^*$  in `TestAllSwaddlings`.  $\tilde{\mathcal{C}}$  intercepts all these swaddlings and decrypts them together with all swaddlings held by  $\mathcal{B}^*$ .  $\tilde{\mathcal{C}}$  then recovers the output state.  $\tilde{\mathcal{C}}$  undoes the quantum operation  $\mathcal{F}$  in order to recover  $\tilde{\mathcal{A}}$ 's effective input state before querying  $\mathcal{F}_1$  with the effective input state. This operation is performed using  $U_{\mathcal{F}}^\dagger \in \text{U}(\mathcal{A}_{\text{out}} \otimes \mathcal{A}_a \otimes \mathcal{B}_{\text{out}} \otimes \mathcal{B}_a)$ , where  $U_{\mathcal{F}}$  is the unitary implemented by the circuit as defined in Sect. 5.1. Obviously,  $U_{\mathcal{F}}^\dagger$  has the same circuit complexity than  $U_{\mathcal{F}}$ .

Then,  $\tilde{\mathcal{C}}$  swaddles back the answer to the query using the same keys. The swaddlings are finally sent to  $\tilde{\mathcal{A}}$  before resuming the interaction between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$ . Two things can happen during the execution of `OpenAllSwaddlings`: 1)  $\tilde{\mathcal{A}}$ 's behavior makes the execution **abort**, and 2)  $\tilde{\mathcal{A}}$  and  $\mathcal{B}$  reach the end of the execution in normal conditions. In the first case,  $\tilde{\mathcal{A}}$  may even get the output state of the computation while preventing  $\mathcal{B}$  from recovering his own. When the output state is available to  $\tilde{\mathcal{A}}$ ,  $\tilde{\mathcal{C}}$  will need to call  $\mathcal{F}_1$  (with  $\tilde{\mathcal{A}}$ 's effective input) and  $\mathcal{F}_2$  where the later call is without fairness (i.e.,  $f = 0$ ) when  $\tilde{\mathcal{A}}$  prevents  $\mathcal{B}$  from decrypting his output logical wires.

More precisely, the simulator  $\tilde{\mathcal{C}}$  will proceed as follows for the rest of the simulation:

3. When  $\tilde{\mathcal{C}}$  receives from  $\mathcal{B}^*$  all of  $\tilde{\mathcal{A}}$ 's swaddlings in `TestAllSwaddlings`, it decrypts them together with all of  $\mathcal{B}^*$ 's swaddlings. This can be done since  $\tilde{\mathcal{C}}$  knows all (classical) encryption keys. The ancilla held in register  $\mathcal{Z}$  is also decrypted.
4. The output of the computation (i.e., between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$ ) is now available to  $\tilde{\mathcal{C}}$ .  $\tilde{\mathcal{C}}$  then runs  $U_{\mathcal{F}}^\dagger$  upon the output registers in order to recover  $\tilde{\mathcal{A}}$ 's effective input state.
5.  $\tilde{\mathcal{C}}$  queries  $\mathcal{F}_1$  with  $\tilde{\mathcal{A}}$ 's effective input state.  $\tilde{\mathcal{C}}$  then re-swaddles all of  $\tilde{\mathcal{A}}$ 's logical wires (now storing the real-world output of the computation) using all secret keys used by  $\tilde{\mathcal{C}}$  at Step 3.  $\tilde{\mathcal{C}}$  then sends all of  $\tilde{\mathcal{A}}$ 's swaddlings back to  $\tilde{\mathcal{A}}$ .
6.  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$  resume their execution from turn  $s^*$  onward. If  $\tilde{\mathcal{A}}$  **aborts** the execution before  $\mathcal{B}^*$  gets his output state then  $\tilde{\mathcal{C}}$  queries  $\mathcal{F}_2$  with  $f = 0$  otherwise  $\tilde{\mathcal{C}}$  queries  $\mathcal{F}_2$  with  $f = 1$ .  $\tilde{\mathcal{C}}$  deletes all its internal registers except the ones held by  $\tilde{\mathcal{A}}$ .

Let  $\overline{\text{E-ABORT}}$  be the event of not having  $\text{E-ABORT}_s$  at any turn  $s < s^*$ . Let  $[\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \overline{\text{E-ABORT}}]^{\mathcal{O}'}$  ( $\rho_{\text{in}}$ ) be denoting an execution between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}$  upon joint input state  $\rho_{\text{in}}$  that does not abort before turn  $s^*$ . Let  $[\tilde{\mathcal{C}} \otimes \mathcal{D} \wedge \overline{\text{E-ABORT}}]^{\mathcal{F}}$  ( $\rho_{\text{in}}$ ) be denoting an execution between  $\tilde{\mathcal{C}}$  and  $\mathcal{B}$ , where the simulated  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$  do not abort before turn  $s^*$  upon joint input state  $\text{tr}_{\mathcal{B}_{\text{in}}}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}}) \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}}$ .

Next lemma establishes the active security when no early aborting occurs. The following lemma follows from the forcingness of the protocol established in Lemma 5.5.

**Lemma 6.2 (No early abort).** *For any quantum adversary  $\tilde{\mathcal{A}}$  in hybrid protocol  $\hat{\Pi}_{\mathcal{F}}^{\mathcal{O}'} = (\mathcal{A}, \mathcal{B}, \mathcal{O}', m)$  for  $\mathcal{F} : \mathcal{L}(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}) \rightarrow \mathcal{L}(\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}})$ , there exists an adversary  $\tilde{\mathcal{C}}$  in  $\Gamma_{\mathcal{F}}$  (polysize in the size of  $\tilde{\mathcal{A}}$  and  $\mathcal{B}$ ) such that for any  $\rho_{\text{in}} \in \mathcal{D}(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}} \otimes \mathcal{R})$ ,*

$$\Delta \left( [\tilde{\mathcal{C}} \otimes \mathcal{D} \wedge \overline{\text{E-ABORT}}]^{\mathcal{F}} (\rho_{\text{in}}), [\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \overline{\text{E-ABORT}}]^{\mathcal{O}'} (\rho_{\text{in}}) \right) \leq \text{negl}(n) .$$

By symmetry, the same is also true with respect to adversaries  $\tilde{\mathcal{B}}$  in  $\hat{\Pi}_{\mathcal{F}}^{\mathcal{O}'}$  and  $\tilde{\mathcal{D}}$  in  $\Gamma_{\mathcal{F}}$ .

*Proof (sketch).* Remember that turn  $s^*$  is where  $\mathcal{B}^*$  is sending  $\tilde{\mathcal{A}}$ 's swaddlings after having tested them with success in `TestAllSwaddlings`. At Step 3,  $\tilde{\mathcal{C}}$  intercepts this quantum transmission from  $\mathcal{B}^*$  to  $\tilde{\mathcal{A}}$ . Let  $[\tilde{\mathcal{C}} \otimes \mathcal{D}]_{\mathcal{S}3}^{\mathcal{F}}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}})$  be denoting the joint state when  $\tilde{\mathcal{C}}$  reaches Step 3 in the simulation above. As we have seen in the proof of Lemma 6.1, we have  $[\tilde{\mathcal{C}} \otimes \mathcal{D}]_{\mathcal{S}3}^{\mathcal{F}}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}}) = [\tilde{\mathcal{A}} \otimes \mathcal{B}]_{s^*}^{\mathcal{O}'}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}})$  since turn  $s^*$  occurs before any player reveals its encryption keys to the other player. Moreover, Lemma 5.5 tells us that

$$\Delta \left( [\tilde{\mathcal{A}} \otimes \mathcal{B}]_{s^*}^{\mathcal{O}'}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}}), \tilde{\mathcal{E}} \circ \mathcal{C}_{s^*} \circ \mathcal{G}_{s^*} \circ \tilde{\mathcal{E}}_0(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}}) \right) \leq \text{negl}(n) . \quad (5)$$

where  $\tilde{\mathcal{E}}_0$  is defined in Definition 5.1,  $\mathcal{G}_{s^*}$  is the *honest* operation so far, and  $\tilde{\mathcal{E}}$  does not act on any of the swaddlings held by  $\tilde{\mathcal{C}}$  since  $\mathcal{B}^*$  just checked them with success (since  $\overline{\text{E-ABORT}}$ ) and are not held by  $\tilde{\mathcal{A}}$  yet. At this point of the protocol,  $\mathcal{G}_{s^*}$  corresponds to the application of  $U_{\mathcal{F}} \in \mathcal{U}(\mathcal{A}_{\text{in}} \otimes \mathcal{A}_a \otimes \mathcal{B}_{\text{in}} \otimes \mathcal{B}_a)$  where  $\mathcal{A}_a$  and  $\mathcal{B}_a$  are the spaces of the ancilla needed to implement  $\mathcal{F}$  unitarily. Notice that  $\mathcal{C}_{s^*}$  is the authentication encoding using all keys (and not only  $\mathcal{B}^*$ 's keys)

since all  $\tilde{\mathcal{A}}$ 's swaddlings have been tested with success by  $\mathcal{B}^*$  before they were returned to  $\tilde{\mathcal{A}}$ . We have:

$$\tilde{\mathcal{E}} \circ \mathbf{C}_{s^*} \circ \mathcal{G}_{s^*} \circ \tilde{\mathcal{E}}_0(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*}) = \tilde{\mathcal{E}} \circ \mathbf{C}_{s^*} \circ U_{\mathcal{F}} \left( \tilde{\mathcal{E}}_0(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}}) \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*} \otimes |0\rangle\langle 0|^{\mathcal{Z}} \right) U_{\mathcal{F}}^\dagger .$$

Using this in (5) before applying  $U_{\mathcal{F}}^\dagger \circ \mathbf{C}_{s^*}^\dagger$ , as  $\tilde{\mathcal{E}}$  does at Steps 3 and 4, results in:

$$\Delta \left( U_{\mathcal{F}}^\dagger \circ \mathbf{C}_{s^*}^\dagger \circ [\tilde{\mathcal{A}} \otimes \mathcal{B}]_{s^*}^{\mathcal{O}'_1}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*}), \right. \\ \left. \tilde{\mathcal{E}} \circ \left( \tilde{\mathcal{E}}_0(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{R}}) \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*} \otimes |0\rangle\langle 0|^{\mathcal{Z}} \right) \right) \leq \text{negl}(n) .$$

At Step 5,  $\tilde{\mathcal{E}}$  calls the ideal functionality  $\mathcal{F}_1$  evaluating  $\mathcal{F}$  upon  $\tilde{\mathcal{A}}$ 's effective input together with  $\mathcal{B}$ 's input while keeping  $\mathcal{B}$ 's output state:

$$\Delta \left( \mathcal{F}_1 \circ U_{\mathcal{F}}^\dagger \circ \mathbf{C}_{s^*}^\dagger \circ [\tilde{\mathcal{A}} \otimes \mathcal{B}]_{s^*}^{\mathcal{O}'_1}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*}), \right. \\ \left. \tilde{\mathcal{E}} \circ \left( \mathcal{F}_1(\tilde{\mathcal{E}}_0(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}})) \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*} \otimes |0\rangle\langle 0|^{\mathcal{Z}} \right) \right) \leq \text{negl}(n) . \quad (6)$$

At this point,  $\tilde{\mathcal{E}}$  encrypts/authenticates, using the same keys as before, all of  $\tilde{\mathcal{A}}$ 's output wires (including her part of  $\mathcal{Z}$  and  $\mathcal{B}_{\text{in}}^*$ ) held by  $\tilde{\mathcal{E}}$ :

$$\Delta \left( \mathbf{C}_{s^*} \circ \mathcal{F}_1 \circ U_{\mathcal{F}}^\dagger \circ \mathbf{C}_{s^*}^\dagger \circ [\tilde{\mathcal{A}} \otimes \mathcal{B}]_{s^*}^{\mathcal{O}'_1}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*}), \right. \\ \left. \tilde{\mathcal{E}} \circ \mathbf{C}_{s^*} \left( \mathcal{F}_1(\tilde{\mathcal{E}}_0(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}})) \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*} \otimes |0\rangle\langle 0|^{\mathcal{Z}} \right) \right) \leq \text{negl}(n) . \quad (7)$$

Applying Lemma 5.5 once again, results in:

$$\Delta \left( [\tilde{\mathcal{A}} \otimes \mathcal{B}]_{s^*}^{\mathcal{O}'_1}(\rho_{\text{in}}), \tilde{\mathcal{E}} \circ \mathbf{C}_{s^*} \circ \mathcal{F} \left( \tilde{\mathcal{E}}_0(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}}) \right) \right) \leq \text{negl}(n) . \quad (8)$$

Let  $[\tilde{\mathcal{E}} \otimes \mathcal{D}]_{\mathcal{S}6}^{\mathcal{F}}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}}) := \mathbf{C}_{s^*} \circ \mathcal{F}_1 \circ U_{\mathcal{F}}^\dagger \circ \mathbf{C}_{s^*}^\dagger \circ [\tilde{\mathcal{A}} \otimes \mathcal{B}]_{s^*}^{\mathcal{O}'_1}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}} \otimes |0\rangle\langle 0|^{\mathcal{B}_{\text{in}}^*})$  be denoting the joint state when  $\tilde{\mathcal{E}}$  reaches Step 6 in the simulation above. Using the triangle inequality for the trace-norm distance using (7) and (8) results in:

$$\Delta \left( \text{tr}_{\mathcal{B}^*} \left( [\tilde{\mathcal{E}} \otimes \mathcal{D}]_{\mathcal{S}6}^{\mathcal{F}}(\rho_{\text{in}}^{\mathcal{A}_{\text{in}}\mathcal{B}_{\text{in}}\mathcal{R}}) \right), [\tilde{\mathcal{A}} \otimes \mathcal{B}]_{s^*}^{\mathcal{O}'_1}(\rho_{\text{in}}) \right) \leq \text{negl}(n) . \quad (9)$$

Equation (9) implies the statement to prove since the interaction between  $\tilde{\mathcal{A}}$  and  $\mathcal{B}^*$  is resumed once  $\tilde{\mathcal{E}}$  transmitted  $\tilde{\mathcal{A}}$ 's swaddlings back to  $\tilde{\mathcal{A}}$ . It follows that  $\tilde{\mathcal{A}}$ 's further actions will produce an output at negligible distance from  $[\tilde{\mathcal{A}} \otimes \mathcal{B} \wedge \overline{\text{E-ABORT}}]_{s^*}^{\mathcal{O}'_1}(\rho_{\text{in}})$  and since  $\tilde{\mathcal{E}}$  deletes all of  $\mathcal{B}^*$ 's registers before getting its output.

To prove the statement for adversary  $\tilde{\mathcal{B}}$  interacting with  $\mathcal{A}$  in  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{O}'_1}$  is easier than for adversary  $\tilde{\mathcal{A}}$ . Remember that for  $\tilde{\mathcal{D}}$ ,  $s^*$  is defined to be the turn at which  $\tilde{\mathcal{B}}$  receives all his swaddlings from  $\mathcal{A}$ . In this case,  $\tilde{\mathcal{D}}$  always allows  $\mathcal{A}$  to get the final state of the computation since the protocol has no lack of fairness for  $\mathcal{A}$  (i.e.  $f = 1$  always). Apart from this difference,  $\tilde{\mathcal{D}}$  proceeds in a very similar way than  $\tilde{\mathcal{E}}$ .  $\square$

## 7 Conclusion

Since Lemmas 6.1 and 6.2 together show that a simulator succeeds in reproducing any real execution, we conclude the active security of the protocol:

**Theorem 7.1 (Active security).** *For any polynomial-time quantum operation  $\mathcal{F} : L(\mathcal{A}_{\text{in}} \otimes \mathcal{B}_{\text{in}}) \rightarrow L(\mathcal{A}_{\text{out}} \otimes \mathcal{B}_{\text{out}})$ , the two-party hybrid protocol  $\widehat{\Pi}_{\mathcal{F}}^{\mathcal{E}'}$  is statistically active secure without fairness for Bob.*

Since the simulator establishing Theorem 7.1 is poly-time, the result in [HSS11] allows us to conclude that instantiating all TCP oracles with computationally secure classical protocols preserves security of our protocol, provided the computational assumption is secure against quantum adversaries.

Our construction can be seen as a *compiler* that transforms protocols secure against specious adversaries into protocols for the same task secure against any adversary, in a spirit somewhat similar to Section 7.4 in [Gol04]. The technique we used allows us to force all players to behave speciously. It would be interesting to find out how general is this compiler in the computational setting.

## References

- ABE08. D. Aharonov, M. Ben-Or, and E. Eban. Interactive proofs for quantum computations. In *Proceedings of Innovations in Computer Science*, 2008. [arxiv.org/abs/0810.5375](https://arxiv.org/abs/0810.5375).
- BCG<sup>+</sup>02. Howard Barnum, Claude Crépeau, Daniel Gottesman, Adam Smith, and Alain Tapp. Authentication of quantum messages. In *43rd Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 449–458, 2002.
- BK05. Sergei Bravyi and Alexei Kitaev. Universal quantum computation with ideal clifford gates and noisy ancillas. *Physical Review A*, 71(022316), 2005. [quant-ph/0403025](https://arxiv.org/abs/quant-ph/0403025).
- BM02. Michael Ben-Or and Dominic Mayers. Quantum universal composability, November 2002. Presentation at "Quantum Information and Cryptography" Workshop, slides online available at <http://www.msri.org/publications/ln/msri/2002/quantumcrypto/mayers/1/meta/aux/mayers.pdf>.
- Can00. Ran Canetti. Security and composition of multiparty cryptographic protocols. *Journal of Cryptology*, 13(1):143–202, 2000.
- CKR09. Matthias Christandl, Robert König, and Renato Renner. Postselection technique for quantum channels with applications to quantum cryptography. *Physical Review Letters*, 102:020504, January 2009.
- Csi98. Imre Csiszár. The method of types. *IEEE Transactions on Information Theory*, 44(6):2505–2523, October 1998.
- DNS10. Frédéric Dupuis, Jesper Buus Nielsen, and Louis Salvail. Secure two-party quantum evaluation of unitaries against specious adversaries. In Tal Rabin, editor, *CRYPTO*, volume 6223 of *Lecture Notes in Computer Science*, pages 685–706. Springer, 2010.
- GC99a. Daniel Gottesman and Isaac L. Chuang. Demonstrating the viability of universal quantum computation using teleportation and single-qubit operations. *Nature*, 402:390–393, November 1999.
- GC99b. Daniel Gottesman and Isaac L. Chuang. Quantum teleportation is a universal computational primitive. [arxiv.org/abs/quant-ph/9908010](https://arxiv.org/abs/quant-ph/9908010), August 1999.
- Gol04. Oded Goldreich. *Foundations of Cryptography*, volume II: Basic Applications. Cambridge University Press, 2004.
- Got97. Daniel Gottesman. *Stabilizer codes and quantum error correction*. PhD thesis, California Institute of Technology, 1997.
- Got10. Daniel Gottesman. An introduction to quantum error correction and fault-tolerant quantum computation. In Jr. Samuel J. Lomonaco, editor, *Quantum Information Science and Its Contributions to Mathematics*, volume 68, pages 13–60. Proceedings of Symposia in Applied Mathematics, April 2010. [arxiv.org/abs/0904.2557](https://arxiv.org/abs/0904.2557).

- GW07. G. Gutoski and J. Watrous. Toward a general theory of quantum games. In *39th Annual ACM Symposium on Theory of Computing (STOC)*, pages 565–574, 2007.
- HSS11. Sean Hallgren, Adam Smith, and Fang Song. Classical cryptographic protocols in a quantum world. In Phillip Rogaway, editor, *CRYPTO*, volume 6841 of *Lecture Notes in Computer Science*, pages 411–428. Springer, 2011.
- LN11. Carolin Lunemann and Jesper Buus Nielsen. Fully simulatable quantum-secure coin-flipping and applications. In Abderrahmane Nitaj and David Pointcheval, editors, *AFRICACRYPT*, volume 6737 of *Lecture Notes in Computer Science*, pages 21–40. Springer, 2011.
- Reg05. Oded Regev. On lattices, learning with errors, random linear codes, and cryptography. In *37th Annual ACM Symposium on Theory of Computing (STOC)*, pages 84–93, 2005.
- Sho96. Peter W. Shor. Fault-tolerant quantum computation. In *37th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 56–65, 1996.
- Unr10. Dominique Unruh. Universally composable quantum multi-party computation. In *Advances in Cryptology—EUROCRYPT 2010*, volume 6110 of *Lecture Notes in Computer Science*, pages 486–505. Springer, 2010.

## A Two-Party Strategies

We define two-party strategies the same way than in [DNS10,GW07]. These strategies may have access to oracles implementing some functionalities, like a communication channel between parties or any other more evolved and non-trivial two-party functionality. The only difference with [DNS10] is that we explicitly allow for  $\mathcal{A}$  and  $\mathcal{B}$  to apply one final operation after the last oracle call. More precisely,

**Definition A.1.** A  $m$ -turn two party strategy with oracle calls denoted  $\Pi^\mathcal{O} = (\mathcal{A}, \mathcal{B}, \mathcal{O}, m)$  consists of:

1. input spaces  $\mathcal{A}_0$  and  $\mathcal{B}_0$  for parties  $\mathcal{A}$  and  $\mathcal{B}$  respectively, together with memory spaces  $\mathcal{A}_1, \dots, \mathcal{A}_{m+1}$  and  $\mathcal{B}_1, \dots, \mathcal{B}_{m+1}$  for  $\mathcal{A}$  and  $\mathcal{B}$  respectively. Spaces  $\mathcal{A}_i$  and  $\mathcal{B}_i$  can be written as  $\mathcal{A}_i = \mathcal{A}_i^\mathcal{O} \otimes \mathcal{A}'_i$  and  $\mathcal{B}_i = \mathcal{B}_i^\mathcal{O} \otimes \mathcal{B}'_i$ , ( $1 \leq i \leq m$ ), and  $\mathcal{O} = (\mathcal{O}_1, \mathcal{O}_2, \dots, \mathcal{O}_m)$  is an  $m$ -tuple of quantum operations:  $\mathcal{O}_i : \mathbb{L}(\mathcal{A}_i^\mathcal{O} \otimes \mathcal{B}_i^\mathcal{O}) \mapsto \mathbb{L}(\mathcal{A}_i^\mathcal{O} \otimes \mathcal{B}_i^\mathcal{O})$ , ( $1 \leq i \leq m$ ).
2. an  $m + 1$ -tuple of quantum operations  $(\mathcal{A}_1, \dots, \mathcal{A}_{m+1})$  for  $\mathcal{A}$ ,  $\mathcal{A}_i : \mathbb{L}(\mathcal{A}_{i-1}) \mapsto \mathbb{L}(\mathcal{A}_i)$ , ( $1 \leq i \leq m + 1$ ),
3. an  $m + 1$ -tuple of quantum operations  $(\mathcal{B}_1, \dots, \mathcal{B}_{m+1})$  for  $\mathcal{B}$ ,  $\mathcal{B}_i : \mathbb{L}(\mathcal{B}_{i-1}) \mapsto \mathbb{L}(\mathcal{B}_i)$ , ( $1 \leq i \leq m + 1$ ).

## B Discussion of our Security Definition

In extending the model in [DNS10] the following changes have been made.

1. The adversary is no more restricted in how it deviates from the protocol, i.e., it can deviate from the protocol as it wants, except that it has to run in polynomial time.
2. The simulator is now allowed to do an input substitution. This is needed as an active attacker might replace its original input by some alternative input and then run the input honestly. This will be indistinguishable from an honest party running with the alternative input, hence the protocol will and must complete successfully. To simulate this “attack” the simulator must be able to ask the ideal functionality to use the alternative input on behalf of corrupted party.

3. In [DNS10], it was required that each step of the protocol could be simulated. This is needed with passive security, as the fact that the view at the end can be simulated is not the same as the view at an earlier point can be simulated, due to no-cloning (cf., [DNS10]). However, an active adversary is allowed to stop its attack at any intermediary step and output its state at that step. Hence security at the end against all adversary implies security at all intermediary steps against all adversaries.
4. When we compare the outputs of the simulation and the protocol, we will compare the entire view, not just the view of the corrupted party. This is usual in defining active security, as it models privacy *and* correctness. In [DNS10] only the view of the adversary and the simulator was compared, which is sufficient to model privacy, as it shows that whatever the adversary learns can be computed given what it is allowed to learn. When comparing also the output of the honest party in the two settings, we also show that whatever influence the adversary can have on the output of the honest party in the protocol, it could also have on the output of the honest party in the simulation. This captures that basically only input substitution is possible, which is the definition of correctness. Furthermore, privacy and correctness cannot meaningfully be separated into separate notions, which is why we capture them using one definition (cf., the discussion on page 150 in [Can00]).
5. Finally, in [DNS10] a simulator contained an explicit bit  $q$  which specified whether it would like to call the ideal functionality or not. This was needed to allow the simulator to preserve its original input. An active adversary has the power to just run the ideal functionality on a dummy input and move its original input to an ancilla if it so desires, so we can remove the explicit bit  $q$ .

Our model is reminiscent of previous models of quantum secure two-party computation, but as opposed to the UC models from [BM02,Unr10,HSS11] our model only guarantees sequential security. That said, since our simulators do not access the code of the adversary or rewind the adversary, our protocol might very well be UC secure in one of the models from [BM02,Unr10,HSS11]. However, here we want to focus on what we consider the main novelty, namely that actively secure 2-party quantum computation is possible at all, and how, and we prefer to avoid cluttering that presentation by dealing with the connoisseur details needed to obtain and prove UC security. We think this is best done in the above simple model.

## C Security of the Clifford-Based QAS

For completeness, we provide the proof of security of the Clifford-based authentication scheme of [ABE08]:

**Theorem C.1 (Security of Clifford-based QAS).** *The Clifford-based QAS defined in Definition 2.3 is  $\varepsilon$ -secure for  $\varepsilon = 6 \times 2^{-n}$ .*

*Proof.* Let  $U_{CR}$  be the attack unitary, and let us decompose it into the Pauli basis on  $C$ :

$$U_{CR} = \sum_P P_C \otimes U_R^P .$$

We also define the following:

$$\mathcal{U}_{R \rightarrow R}^{\text{acc}}(\rho) = U_R^{\mathbb{1}} \rho U_R^{\mathbb{1}\dagger}$$

$$\begin{aligned}\mathcal{U}_{R \rightarrow R}^{\text{rej}}(\rho) &= \sum_{P \neq 1} U_R^P \rho U_R^{P\dagger} \\ \Omega_S &= \pi_S .\end{aligned}$$

We now compute the state given that the decoder accepts:

$$\begin{aligned}& P_{\text{acc}} \mathbb{E}_k \mathcal{D}_k \left( U_{CR} \mathcal{E}_k(\psi_{SR}) U_{CR}^\dagger \right) P_{\text{acc}} \\ &= \mathbb{E}_k \sum_{P, P'} \text{tr}_A \left[ P_{\text{acc}} C_k^\dagger(P_C \otimes U_R^P) C_k(\psi_{SR} \otimes |0^n\rangle\langle 0^n|_A) C_k^\dagger(P_C^\dagger \otimes U_R^{P'\dagger}) C_k P_{\text{acc}} \right] \\ &\stackrel{(a)}{=} \mathbb{E}_k \sum_P \text{tr}_A \left[ P_{\text{acc}} C_k^\dagger(P_C \otimes U_R^P) C_k(\psi_{SR} \otimes |0^n\rangle\langle 0^n|_A) C_k^\dagger(P_C^\dagger \otimes U_R^{P\dagger}) C_k P_{\text{acc}} \right] \\ &= U_R^1 \psi_{SR} U_R^{1\dagger} + \mathbb{E}_k \sum_{P \neq 1} \text{tr}_A \left[ P_{\text{acc}} C_k^\dagger(P_C \otimes U_R^P) C_k(\psi_{SR} \otimes |0^n\rangle\langle 0^n|_A) C_k^\dagger(P_C^\dagger \otimes U_R^{P\dagger}) C_k P_{\text{acc}} \right] \\ &= \mathcal{U}^{\text{acc}}(\psi_{SR}) + \sum_{P \neq 1} \mathbb{E}_{\tilde{P} \neq 1} \text{tr}_A \left[ P_{\text{acc}}(\tilde{P}_C \otimes U_R^P)(\psi_{SR} \otimes |0^n\rangle\langle 0^n|_A)(\tilde{P}_C^\dagger \otimes U_R^{P\dagger}) P_{\text{acc}} \right] \\ &= \mathcal{U}^{\text{acc}}(\psi_{SR}) + \frac{2^{-n} |\mathfrak{P}|}{|\mathfrak{P}| - 1} \sum_{P \neq 1} U_R^P (\pi_S \otimes \psi_R) U_R^{P\dagger} - \frac{1}{|\mathfrak{P}| - 1} \sum_{P \neq 1} U_R^P \psi_{SR} U_R^{P\dagger} \\ &= \mathcal{U}^{\text{acc}}(\psi_{SR}) + \frac{2^{-a} |\mathfrak{P}|}{|\mathfrak{P}| - 1} \mathcal{U}^{\text{rej}}(\psi_R) \otimes \pi_S - \frac{1}{|\mathfrak{P}| - 1} \mathcal{U}^{\text{rej}}(\psi_{SR}) ,\end{aligned}$$

where (a) comes from Lemma B.3 in [ABE08], and where  $\mathfrak{P}$  denotes the set of all Pauli operators on  $n$  qubits. Likewise, the state given that the decoder rejects is given by:

$$\begin{aligned}& P_{\text{rej}} \mathbb{E}_k \mathcal{D}_k \left( U_{CR} \mathcal{E}_k(\psi_{SR}) U_{CR}^\dagger \right) P_{\text{rej}} \\ &= \mathbb{E}_k \sum_{P, P'} \text{tr}_A \left[ P_{\text{rej}} C_k^\dagger(P_C \otimes U_R^P) C_k(\psi_{SR} \otimes |0^n\rangle\langle 0^n|_A) C_k^\dagger(P_C^\dagger \otimes U_R^{P'\dagger}) C_k P_{\text{rej}} \right] \\ &\stackrel{(a)}{=} \mathbb{E}_k \sum_P \text{tr}_A \left[ P_{\text{rej}} C_k^\dagger(P_C \otimes U_R^P) C_k(\psi_{SR} \otimes |0^n\rangle\langle 0^n|_A) C_k^\dagger(P_C^\dagger \otimes U_R^{P\dagger}) C_k P_{\text{rej}} \right] \\ &= \mathbb{E}_k \sum_{P \neq 1} \text{tr}_A \left[ P_{\text{rej}} C_k^\dagger(P_C \otimes U_R^P) C_k(\psi_{SR} \otimes |0^n\rangle\langle 0^n|_A) C_k^\dagger(P_C^\dagger \otimes U_R^{P\dagger}) C_k P_{\text{rej}} \right] \\ &= \sum_{P \neq 1} \mathbb{E}_{\tilde{P} \neq 1} \text{tr}_A \left[ P_{\text{rej}}(\tilde{P}_C \otimes U_R^P)(\psi_{SR} \otimes |0^n\rangle\langle 0^n|_A)(\tilde{P}_C^\dagger \otimes U_R^{P\dagger}) P_{\text{rej}} \right] \\ &= \frac{(1 - 2^{-n}) |\mathfrak{P}|}{|\mathfrak{P}| - 1} \sum_{P \neq 1} U_R^P (\pi_S \otimes \psi_R) U_R^{P\dagger} - \frac{1}{|\mathfrak{P}| - 1} \sum_{P \neq 1} U_R^P \psi_{SR} U_R^{P\dagger} \\ &= \frac{(1 - 2^{-n}) |\mathfrak{P}|}{|\mathfrak{P}| - 1} \mathcal{U}^{\text{rej}}(\psi_R) \otimes \pi_S - \frac{1}{|\mathfrak{P}| - 1} \mathcal{U}^{\text{rej}}(\psi_{SR}) \\ &= \mathcal{U}^{\text{rej}}(\psi_R) \otimes \pi_S - \frac{2^{-a} |\mathfrak{P}| - 1}{|\mathfrak{P}| - 1} \mathcal{U}^{\text{rej}}(\psi_R) \otimes \pi_S - \frac{1}{|\mathfrak{P}| - 1} \mathcal{U}^{\text{rej}}(\psi_{SR}) .\end{aligned}$$



Now, it is easy to see that the trace distance in Definition 2.2 can be upper-bounded by:

$$\begin{aligned}
\text{tr} \left[ 2 \times \frac{2^{-n} |\mathfrak{P}|}{|\mathfrak{P}| - 1} \mathcal{U}^{\text{rej}}(\psi_R) \otimes \pi_S + \frac{2}{|\mathfrak{P}| - 1} \mathcal{U}^{\text{rej}}(\psi_{SR}) \right] &\leq 2 \times \frac{2^{-n} |\mathfrak{P}|}{|\mathfrak{P}| - 1} + \frac{2}{|\mathfrak{P}| - 1} \\
&= 2 \times \frac{2^{-n} 4^{n+s}}{4^{n+s} - 1} + \frac{2}{4^{n+s} - 1} \\
&\leq 6 \times 2^{-n} .
\end{aligned}$$

□

## D R-gate teleportation and magic state distillation

The gate teleportation protocol for the R-gate works as follows. Suppose one wants to execute R on wire  $w$ . One would need to initialize a wire  $m$  in the *magic state*  $|M\rangle := \frac{1}{\sqrt{2}} (|0\rangle + e^{i\pi/4}|1\rangle)$ . The protocol then consists of a CNOT from  $m$  to  $w$ , and then of a measurement of  $w$  in the computational basis. If the result is 0, one does nothing, and if the result is 1, we apply the gate  $A := e^{-i\pi/4}PX$  to  $m$ . We then rename  $m$  to  $w$ , and one can check that an R-gate has been applied to  $w$ .

We therefore need a supply of magic states for our protocol. However, we need some way of generating them securely. For this, we turn to ideas from fault-tolerant computation, more specifically from [BK05]. There, they consider a model of computation in which we can perform perfect Clifford gates, but where we only have a supply of noisy magic states. We therefore need to distill good magic states from noisy ones. We will use this technique to ensure that we have a supply of good magic states: we will first let Alice prepare a large supply of magic states, and then we will have Bob distill a smaller number of states from this untrusted supply. Bob will first test a random sample of the states and check that there are no errors in the sample, and he will then use the distillation protocol from [BK05] to get the final states.

Here we will simply state the protocol from [BK05] without attempting to explain how and why it works. The protocol is based on a CSS code which is constructed from two linear classical codes  $\mathcal{L}_1$  and  $\mathcal{L}_2$ , whose generator matrices are given in Figure D.1 along with the stabilizer generators of the resulting CSS codes.

Here's how one step of the magic state distillation works:

1. Dephase all qubits (apply  $A$  to each qubit with probability 1/2).
2. Randomly permute all the qubits.
3. Measure  $\mu$  (the syndrome of  $\mathcal{L}_2$ ).
4. Compute the vector  $w$  as  $w = \sum_g \mu(g)g$ , where the sum ranges over stabilizer generators of  $\mathcal{L}_2$ .
5. Apply  $A(w)^\dagger$ .
6. Measure  $\eta$  (syndrome of  $\mathcal{L}_1$  with  $X$ 's)
7. If  $\eta \neq 0$ , declare an error.
8. Apply the decoding operation for the  $X/Z$  CSS code.
9. Dephase all the qubits again.

In [BK05], the protocol is analyzed assuming that the input consists of  $n$  independent, identically distributed states of the form  $\rho = (1 - \varepsilon)|M\rangle\langle M| + \varepsilon|\bar{M}\rangle\langle\bar{M}|$ , where  $|\bar{M}\rangle := \frac{1}{\sqrt{2}}(|0\rangle -$

$$\begin{bmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

(a) Generator matrix of code  $\mathcal{L}_1$

$$\begin{bmatrix} 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 0 & 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

(b) Generator matrix of code  $\mathcal{L}_2$

$$\begin{bmatrix} X & I & X & I & X & I & X & I & X & I & X & I & X & I & X \\ I & X & X & I & I & X & X & I & I & X & X & I & I & X & X \\ I & I & I & X & X & X & X & I & I & I & I & X & X & X & X \\ I & I & I & I & I & I & I & X & X & X & X & X & X & X & X \\ Z & I & Z & I & Z & I & Z & I & Z & I & Z & I & Z & I & Z \\ I & Z & Z & I & I & Z & Z & I & I & Z & Z & I & I & Z & Z \\ I & I & I & Z & Z & Z & Z & I & I & I & I & Z & Z & Z & Z \\ I & I & I & I & I & I & I & Z & Z & Z & Z & Z & Z & Z & Z \\ I & I & Z & I & I & I & Z & I & I & I & Z & I & I & I & Z \\ I & I & I & I & Z & I & Z & I & I & I & I & I & Z & I & Z \\ I & I & I & I & I & I & I & I & Z & I & Z & I & Z & I & Z \\ I & I & I & I & I & Z & Z & I & I & I & I & I & I & Z & Z \\ I & I & I & I & I & I & I & I & I & Z & Z & I & I & Z & Z \\ I & I & I & I & I & I & I & I & I & I & Z & Z & Z & Z & Z \end{bmatrix}$$

(c) Stabilizer generators of the CSS code used for magic state distillation

Fig. D.1: Description of the code used for magic state distillation.

$e^{i\pi/4}|1\rangle\rangle$  is a state orthogonal to  $|M\rangle$ . If we assume this, then it can then be shown that, if the input has an error rate of  $\varepsilon$ , then the output error rate is  $\varepsilon_{\text{out}} \leq 40\varepsilon^3$  whenever  $\varepsilon \leq 0.041$  (see [BK05], Equation (36)). This means that after repeating this recursively  $\ell$  times, we have a final error rate of  $\varepsilon_f \leq \frac{1}{\sqrt{40}} (\sqrt{40\varepsilon})^{3^\ell}$ .

However, our scenario is somewhat different from this: we do not have an source of independently and identically distributed magic states with a predefined error rate. Instead, we control the error rate of the magic states by sampling a fraction of them, and reject if some of them are not  $|M\rangle$ . This gives us that with very high probability, the set of magic states after sampling lies in a subspace containing at most some fraction  $\varepsilon$  of errors. That the distillation protocol nevertheless works in this case follows from Theorem D.1 below, which shows that if the protocol has a low probability of failure on iid states with a low error rate, it also has a low probability of failure on states supported on a low-error subspace.

**Theorem D.1.** *Let  $\mathcal{H} = \text{span}\{|0\rangle, |1\rangle\}$ ,  $\mathcal{H}_t := \text{span}\{|w\rangle : w \in \{0, 1\}^n, |w| = t\} \subset \mathcal{H}^{\otimes n}$ , and let  $\mathcal{H}_{\leq t} := \bigoplus_{s \leq t} \mathcal{H}_s$ . Furthermore, let  $\rho \in \text{D}(\mathcal{H}_{\leq t})$  be diagonal in the computational basis, let  $\mathcal{P} : \text{L}(\mathcal{H}^{\otimes n}) \rightarrow \text{L}(\mathcal{H}^{\otimes m})$  be a CPTP map such that  $\mathcal{P}(\pi\omega\pi^\dagger) = \mathcal{P}(\omega)$  for any permutation  $\pi$  of the input qubits and any  $\omega \in \text{Herm}(\mathcal{H}^{\otimes n})$  and whose output is guaranteed to be diagonal in the computational basis. Then,*

$$\|\mathcal{P}(\rho) - |0\rangle\langle 0|^{\otimes m}\|_1 \leq (n+1) \max_{s \leq t} \|\mathcal{P}(\sigma_s^{\otimes n}) - |0\rangle\langle 0|^{\otimes m}\|_1,$$

where  $\sigma_s = \frac{n-s}{n}|0\rangle\langle 0| + \frac{s}{n}|1\rangle\langle 1|$ .

This theorem is somewhat reminiscent of the ‘‘postselection’’ technique from [CKR09], which can be used to reduce the security of QKD against general attacks to security against collective attacks.

*Proof.* Without loss of generality, assume that  $\rho$  is invariant under permutations of the  $n$  qubits (otherwise replace it by  $\mathbb{E}_\pi \pi \rho \pi^\dagger$ ). Then,  $\rho$  is of the form

$$\rho = \sum_{s=0}^t \lambda_s \binom{n}{s}^{-1} P_s \leq (n+1) \sum_{s \leq t} \lambda_s \sigma_s^{\otimes n},$$

where  $P_s$  is the projector onto  $\mathcal{H}_s$ , and  $\{\lambda_s\}$  forms a probability distribution. The inequality above stems from the fact that  $\text{tr}[P_s \sigma_s^{\otimes n}]$  is at least as large as the average of  $\text{tr}[P_{s'} \sigma_s^{\otimes n}]$  over a uniform choice of  $s'$ . (See [Csi98, Lemma II.2] for a proof.) Hence, we now have that:

$$\begin{aligned} \|\mathcal{P}(\rho) - |0\rangle\langle 0|^{\otimes m}\|_1 &= 2 \text{tr}[(\mathbb{1} - |0\rangle\langle 0|^{\otimes m}) \mathcal{P}(\rho)] \\ &\leq 2(n+1) \text{tr} \left[ (\mathbb{1} - |0\rangle\langle 0|^{\otimes m}) \mathcal{P} \left( \sum_{s \leq t} \lambda_s \sigma_s^{\otimes n} \right) \right] \\ &\leq 2(n+1) \max_{s \leq t} \text{tr} [(\mathbb{1} - |0\rangle\langle 0|^{\otimes m}) \mathcal{P}(\sigma_s^{\otimes n})] \\ &= (n+1) \max_{s \leq t} \|\mathcal{P}(\sigma_s^{\otimes n}) - |0\rangle\langle 0|^{\otimes m}\|_1. \end{aligned}$$

In the above, the two equalities are due to the fact that the output of  $\mathcal{P}$  is guaranteed to be diagonal in the computational basis, and that it is trace preserving. This concludes the proof.  $\square$

To apply this theorem to the problem at hand, we let  $\mathcal{P}$  be the distillation protocol after the sampling step and after having dephased all the input states, set  $t = \lfloor 0.041n \rfloor$ , and replace the computational basis with the  $\{|M\rangle, |\bar{M}\rangle\}$  basis. The dephasing and permutation steps ensure that  $\mathcal{P}$  satisfies the requirements of the theorem. We have now shown that if the number of errors in the input is less than or equal to  $0.041n$ , the final error rate will be negligible. All that is now left to do is to ensure that, if the actual error rate was above 0.041, then the adversary would have been caught with overwhelming probability. Note that since the alleged magic states were dephased in the  $A$ -basis, this corresponds to a purely classical sampling problem: if one has a set of  $n$  states, with a fraction of at least 0.041 being faulty, what is the probability  $p$  that sampling  $k$  of them reveals only good states? This can be calculated as follows:

$$p = \binom{\lfloor 0.959n \rfloor}{k} \binom{n}{k}^{-1} \leq \frac{2^{\lfloor 0.959n \rfloor h(k/\lfloor 0.959n \rfloor)}}{\frac{1}{n+1} 2^{nh(k/n)}} \leq (n+1) 2^{0.959nh(k/0.959n) - nh(k/n)}.$$

Here, the inequality follows again from [Csi98, Lemma II.2]. Now, it can be shown that the right-hand side is upper-bounded by  $(n+1)2^{-0.1n}$  whenever  $k/n \geq 0.8$ .

## E Description of the subprotocols

The ancillas swaddled during Initialization also need to have their initial state verified. We now describe the verification procedure `VerifyAncilla` that verifies that a swaddled ancilla is indeed in state  $|0\rangle$ . This subprotocol relies on the ideal functionality  $\text{id}_{\text{INTCNOT}}$  defined as follows:

$$\text{id}_{\text{INTCNOT}}((i_a, i_b), (i'_a, i'_b)) = \begin{cases} (K'_{w,a}, (t, K'_{w,b})) & \text{if } i_a = i'_a = \text{id}_A(w) \text{ and } i_b = i'_b = \text{id}_B(w) \\ (\perp, \perp) & \text{otherwise,} \end{cases}$$

where  $t$  is the description of a Clifford  $T \in \mathfrak{C}_{2n+1}$  such that

$$T = \begin{cases} B'(A' \otimes \mathbb{1}_n)C(A^\dagger \otimes \mathbb{1}_n)B^\dagger & \text{if Alice holds the inner key,} \\ A'(\mathbb{1}_n \otimes B')C(\mathbb{1}_n \otimes B^\dagger)A^\dagger & \text{if Alice holds the outer key,} \end{cases}$$

where  $C$  is a CNOT from the data qubit as control and Bob's first dummy qubit as target, and  $A'$  and  $B'$  are Cliffords of the appropriate size chosen uniformly at random.

**VerifyAncilla**( $\mathfrak{s}(w)$ ), with  $w \in \mathcal{A}$ :

1. Alice sends  $\mathfrak{s}(w)$  to Bob.
2. Bob runs `TestSwaddling`( $\mathfrak{s}(w)$ ).
3. Alice and Bob call  $(K'_{w,a}, (t, K'_{w,b})) \leftarrow \text{id}_{\text{INTCNOT}}((\text{id}_A(w), \text{id}_B(w)), (\text{id}_A(w), \text{id}_B(w)))$ . Bob executes  $T$  on  $\mathfrak{s}(w)$ .
4. Bob sends  $\mathfrak{s}(w)$  back to Alice.
5. Alice runs `TestSwaddling`( $\mathfrak{s}(w)$ ) and sends  $\mathfrak{s}(w)$  to Bob.
6. Bob runs `TestSwaddling`( $\mathfrak{s}(w)$ ).
7. Bob sends  $\mathfrak{s}(w)$  back to Alice.
8. Alice runs `TestSwaddling`( $\mathfrak{s}(w)$ ).

In the above, the idea is to measure the swaddled state in the computational basis by first executing a CNOT with the swaddled state as the control and one of Bob's dummies as the target. Bob then performs the measurement implicitly when he tests his dummies in step 6.

The next procedure implements the magic state distillation outlined in Appendix D.

**DistillMagic( $\mathfrak{s}(\mathfrak{w})$ ), with  $\mathfrak{w} \in \mathcal{A}$ :**

1. Alice creates  $N$  qubits in state  $|M\rangle$ .
2. Swaddle all of them.
3. Apply  $A$  to each qubit with probability  $1/2$ .
4. Measure  $\mu$  (the syndrome of  $\mathcal{L}_2$ ).
5. Compute the vector  $w$  as  $w = \sum_g \mu(g)g$ , where the sum ranges over stabilizer generators of  $\mathcal{L}_2$ .
6. Apply  $A(w)^\dagger$  to the qubits.
7. Measure  $\eta$  (syndrome of  $\mathcal{L}_1$  with  $X$ 's).
8. If  $\eta \neq 0$ , declare an error.
9. Apply the decoding operation for the  $X/Z$  CSS code.
10. Apply  $A$  to each qubit again with probability  $1/2$ .

In the above, the measurement are done by calling the subprotocol Measure, given further down.

**RGate( $\mathfrak{s}(\mathfrak{w})$ ) with  $\mathfrak{w} \in \mathcal{A}$ :**

1. Alice performs a swaddled CNOT with  $\mathfrak{m}_w$  as a control, and  $\mathfrak{w}$  as target.
2. Alice sends  $\mathfrak{s}(\mathfrak{w})$  to Bob. Bob runs **TestSwaddling**( $\mathfrak{s}(\mathfrak{w})$ ) on it.
3. Alice runs **Measure**( $\mathfrak{s}(\mathfrak{w})$ ).
4. Alice and Bob perform a TPC whose result is a Clifford which, if both measurement results were zero, updates the key, and if both measurement results were one, performs a swaddled  $e^{i\pi/4}XP^\dagger$  and then updates the key. If the measurement results differ, then they abort.
5. Alice relabels  $\mathfrak{m}_{\alpha,i}$  to  $\mathfrak{w}$ .

To see that such a procedure indeed performs an R-gate, we will analyze its behavior when the  $\mathfrak{w}$  is in some arbitrary pure state  $\alpha|0\rangle + \beta|1\rangle$ ; hence, if we include the magic state, the initial state here is

$$\frac{1}{\sqrt{2}}(|0\rangle + e^{i\pi/4}|1\rangle)(\alpha|0\rangle + \beta|1\rangle) .$$

After the CNOT, we get

$$\frac{1}{\sqrt{2}}(\alpha|00\rangle + \beta|01\rangle + e^{i\pi/4}\alpha|11\rangle + e^{i\pi/4}\beta|10\rangle) .$$

Measuring the second qubit then yields either  $\alpha|0\rangle + e^{i\pi/4}\beta|1\rangle$  if the measurement result was 0, and  $\beta|0\rangle + e^{i\pi/4}\alpha|1\rangle$  if the result was 1. This is the desired result in the first case, and performing  $e^{i\pi/4}XP^\dagger$  in the second case also yields the same state.

We now know how to perform all of these operations on swaddled states; the only difficulty is to ensure that Alice doesn't lie about the measurement result. Hence, we first copy the measurement result into another swaddled register that Bob gets, ensuring that he also gets the measurement result.

**Measure( $\mathfrak{s}(\mathfrak{w})$ ), with  $\mathfrak{w} \in \mathcal{A}$ :**

1. Let  $\mathfrak{a}$  be an ancilla in state  $|0\rangle$  prepared during **Initialization**, and let Alice call **CNOT**( $\mathfrak{s}(\mathfrak{w})$ ,  $\mathfrak{s}(\mathfrak{a})$ ).
2. Alice sends  $\mathfrak{s}(\mathfrak{a})$  to Bob.
3. Bob runs **TestDummies**( $\mathfrak{s}(\mathfrak{a})$ ).
4. Alice and Bob call  $\text{id}_{\text{OPEN}}(\text{id}_B(\mathfrak{w}), \text{id}_B(\mathfrak{w}))$  and  $\text{id}_{\text{OPEN}}(\text{id}_A(\mathfrak{a}), \text{id}_A(\mathfrak{a}))$ .
5. Alice decrypts  $\mathfrak{w}$  and measures it in the computational basis, and Bob decrypts  $\mathfrak{a}$  and measures it in the computational basis.

This subprotocol measures  $\mathfrak{w}$  in the computational basis and ensures that both Alice and Bob get the result.

## F Additional definitions

**Definition F.1 (Inner swaddling operation).** We define  $C_{I,a,\mathfrak{w}}$ , Alice's inner swaddling operation for wire  $\mathfrak{w}$ , as follows:

$$C_{I,a,\mathfrak{w}}(\rho) := \sum_{C_a \in \mathcal{C}(k+1)} \left( |k_a\rangle\langle k_a|^{\mathcal{K}_a} \otimes C_{k_a}^{\mathfrak{s}(\mathfrak{w})} \right) \rho \left( |k_a\rangle\langle k_a|^{\mathcal{K}_a} \otimes C_{k_a}^{\mathfrak{s}(\mathfrak{w})\dagger} \right),$$

where  $k_a$  is a classical description of  $C_a$ . Likewise,  $C_{I,b,\mathfrak{w}}$  is Bob's inner swaddling operation and is defined the same way.

**Definition F.2 (Outer swaddling operation).** We define  $C_{O,a,\mathfrak{w}}$ , Alice's outer swaddling operation for wire  $\mathfrak{w}$ , as follows:

$$C_{O,a,\mathfrak{w}}(\rho) := \sum_{C_a \in \mathcal{C}(2k+1)} \left( |k_a\rangle\langle k_a|^{\mathcal{K}_a} \otimes C_{k_a}^{\mathfrak{s}(\mathfrak{w})} \right) \rho \left( |k_a\rangle\langle k_a|^{\mathcal{K}_a} \otimes C_{k_a}^{\mathfrak{s}(\mathfrak{w})\dagger} \right),$$

where  $k_a$  is a classical description of  $C_a$ . Likewise,  $C_{O,b,\mathfrak{w}}$  is Bob's outer swaddling operation and is defined the same way.

**Definition F.3 (Swaddling pattern).** For every turn  $s$  of the protocol, let  $\mathfrak{A}_{I,s}$  (resp.  $\mathfrak{B}_{I,s}$ ) be the set of wires for which Alice (resp. Bob) holds the inner key, and  $\mathfrak{A}_{O,s}$  (resp.  $\mathfrak{B}_{O,s}$ ) be the set of wires for which Alice holds the outer key. Wires that are not swaddled yet are in none of these sets, and wires which only have an inner swaddling are only in the corresponding inner set. We call these sets the swaddling pattern at turn  $s$ .

**Definition F.4.** Let  $C_s$  be the operation that takes all the unswaddled wires and swaddles according to the swaddling pattern at turn  $s$ :

$$C_s := \left[ \left( \bigotimes_{\mathfrak{w} \in \mathfrak{A}_{O,s}} C_{O,a,\mathfrak{w}} \right) \otimes \left( \bigotimes_{\mathfrak{w} \in \mathfrak{B}_{O,s}} C_{O,b,\mathfrak{w}} \right) \right] \circ \left[ \left( \bigotimes_{\mathfrak{w} \in \mathfrak{A}_{I,s}} C_{I,a,\mathfrak{w}} \right) \otimes \left( \bigotimes_{\mathfrak{w} \in \mathfrak{B}_{I,s}} C_{I,b,\mathfrak{w}} \right) \right].$$

Furthermore, we define Alice's total swaddling operation:

$$C_{a,s} := \left( \bigotimes_{\mathfrak{w} \in \mathfrak{A}_{O,s}} C_{O,a,\mathfrak{w}} \right) \circ \left( \bigotimes_{\mathfrak{w} \in \mathfrak{A}_{I,s}} C_{I,a,\mathfrak{w}} \right),$$

and likewise for Bob's total swaddling operation  $C_{b,s}$ .

## G Proving forcingness

In the following, we will not give explicit bounds for the various negligible functions, and assume that states that have negligible trace distance to each other are equal. Explicit bounds can be derived from Theorem 2.4, which gives an upper bound on the error of the the Clifford-based QAS.

**Lemma G.1.** Initialization is  $\mathcal{G}$ -forcing for any  $\tilde{\mathcal{A}}$  and any  $\tilde{\mathcal{B}}$ , where  $\mathcal{G}$  tests the ancillas in the  $|0\rangle$  state and distills magic states (see Lemmas G.3 and G.2).

*Proof.* We first prove forcingness against  $\tilde{\mathcal{A}}$ . Let the state just before Alice sends her qubits to Bob in the first call to Swaddle be

$$\rho_1^{ABR} = \tilde{\mathcal{E}}_0(\rho_0^{ABR}),$$

for some operation  $\tilde{\mathcal{E}}_0$  on Alice's systems. At this point, Alice commits to a Clifford  $C_A$ , sends  $n + 1$  qubits to Bob and receives  $\mathfrak{s}(w)$  from Bob. Alice then gets the qubits back from Bob, along with more qubits in the second call to Swaddle. At this point, the global state is

$$\rho_2^{ABR} = C_{b,s} \circ \tilde{\mathcal{E}}_0(\rho_0^{ABR}).$$

She will then send back some of these qubits to Bob. Now, since we are looking at the case where there is no abortion before the end, the fact that she passes Bob's `TestDummies` means that her attack can be represented by a CP map  $\mathcal{U}^{\text{acc}}$  on her qubits that she is not sending, as in Definition 2.2. The state after sending the qubits back to Bob is therefore

$$\rho_3^{ABR} = \mathcal{U}^{\text{acc}} \circ C_{b,s} \circ \tilde{\mathcal{E}}_0(\rho_0^{ABR}).$$

and the protocol is therefore  $\mathbb{1}$ -forcing for  $\tilde{\mathcal{A}}$  so far. The rest of the protocol only involves calls to other subprotocols which are proven to be forcing later on.

The proof of forcingness against  $\tilde{\mathcal{B}}$  follows the same argument.  $\square$

**Lemma G.2.** `VerifyAncilla`( $\mathfrak{s}(w)$ ) is  $\mathcal{G}$ -forcing for any  $\tilde{\mathcal{A}}$  and any  $\tilde{\mathcal{B}}$ , where  $\mathcal{G}(\cdot) = |0\rangle\langle 0| \cdot |0\rangle\langle 0|$ .

*Proof.* We first prove forcingness against  $\tilde{\mathcal{A}}$ . Assume that the state at the beginning of the protocol (in the no-early-abort case) is

$$\rho^{ABR} = \tilde{\mathcal{E}} \circ C_{b,s} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0(\rho_0^{ABR}).$$

Since the state passes Bob's `TestDummies` in step 2, one can assume that  $\tilde{\mathcal{A}}$ 's attack in step 1 can be represented by some  $\mathcal{U}_1^{\text{acc}}$  on her systems. Furthermore, Bob implemented a measurement of the logical qubit in the computational basis, via his own dummy: since his dummy was initialized to  $|0\rangle$  the CNOT that was performed from the logical qubit to the dummy, followed by measuring all his dummies in the computational basis means that the logical qubit was also measured. Since for forcingness, we only consider the case in which all `TestDummies` are passed, it follows that the operation  $\mathcal{G}$  defined above was applied. Again, in step 5, Alice must have sent something that passes Bob's `TestDummies`, so her operation can be represented by some  $\mathcal{U}_2^{\text{acc}}$  on her systems. Finally, in step 8, she can again do an arbitrary operation  $\tilde{\mathcal{E}}_f$  on her systems. Hence, the final state looks like

$$\rho_f^{ABR} = \tilde{\mathcal{E}}' \circ C_{b,s} \circ \mathcal{G} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0(\rho_0^{ABR}),$$

where  $\tilde{\mathcal{E}}' := \tilde{\mathcal{E}}_f \circ \mathcal{U}_2^{\text{acc}} \circ \mathcal{U}_1^{\text{acc}} \circ \tilde{\mathcal{E}}$ .

We now prove forcingness against  $\tilde{\mathcal{B}}$ . Assume that the state at the beginning is

$$\rho^{ABR} = \tilde{\mathcal{E}} \circ C_{a,s} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0(\rho_0^{ABR}).$$

Since Bob passes the `TestDummies` in step 5, his attack must correspond to a  $\mathcal{U}_1^{\text{acc}}$  on his other systems; likewise at step 8 with  $\mathcal{U}_2^{\text{acc}}$ . The final state can therefore be expressed as

$$\rho_f^{ABR} = \tilde{\mathcal{E}}' \circ C_{a,s} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0 (\rho_0^{ABR}),$$

where  $\tilde{\mathcal{E}}' := \mathcal{U}_2^{\text{acc}} \circ \mathcal{U}_1^{\text{acc}} \circ \tilde{\mathcal{E}}$ . Finally, since Alice is honest, she has indeed prepared the ancilla in state  $|0\rangle$ , so adding  $\mathcal{G}$  in the above expression leaves it invariant, and we have

$$\rho_f^{ABR} = \tilde{\mathcal{E}}' \circ C_{a,s} \circ \mathcal{G} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0 (\rho_0^{ABR})$$

as the final states.  $\square$

**Lemma G.3.** `DistillMagic`( $\mathfrak{s}(\mathbf{w})$ ) is  $\mathcal{G}$ -forcing for any  $\tilde{\mathcal{A}}$  and any  $\tilde{\mathcal{B}}$ , where  $\mathcal{G}$  is the magic state distillation protocol described in Appendix D.

*Proof.* The distillation protocol is simply a combination of other subprotocols, so the forcingness of the subprotocols used guarantee the forcingness of the distillation protocol.  $\square$

**Lemma G.4.** `CNOT`( $\mathfrak{s}(\mathbf{w}_c), \mathfrak{s}(\mathbf{w}_t)$ ) with  $\mathbf{w}_c$  at Alice, is  $\mathcal{G}$ -forcing for any  $\tilde{\mathcal{A}}$  and any  $\tilde{\mathcal{B}}$ , where  $\mathcal{G}$  is a CNOT on the appropriate wires.

*Proof.* We start by proving forcingness for  $\tilde{\mathcal{A}}$ . Assume that the state at the beginning of the protocol (in the no-early-abort case) is

$$\rho^{ABR} = \tilde{\mathcal{E}} \circ C_{b,s} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0 (\rho_0^{ABR}).$$

Since the state passes Bob's joint `TestDummies` in step 3, one can assume that  $\tilde{\mathcal{A}}$ 's attack in step 2 can be represented by some  $\mathcal{U}^{\text{acc}}$  on her systems. Bob then applies the CNOT honestly, so it follows that the operation  $\mathcal{G}$  defined above was applied. Finally, in step 8, she can again do an arbitrary operation  $\tilde{\mathcal{E}}_f$  on her systems. Hence, the final state looks like

$$\rho_f^{ABR} = \tilde{\mathcal{E}}' \circ C_{b,s} \circ \mathcal{G} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0 (\rho_0^{ABR}),$$

where  $\tilde{\mathcal{E}}' := \tilde{\mathcal{E}}_f \circ \mathcal{U}^{\text{acc}} \circ \tilde{\mathcal{E}}$ .

We now prove forcingness for  $\tilde{\mathcal{B}}$ ; we will suppose that Bob holds  $\mathbf{w}_t$ , if not one can easily see that the proof is a slightly simpler version of the same argument. Assume that the state at the beginning of the protocol is

$$\rho^{ABR} = \tilde{\mathcal{E}} \circ C_{a,s} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0 (\rho_0^{ABR}).$$

Since the state passes Alice's `TestDummies` in step 1, one can assume that  $\tilde{\mathcal{B}}$ 's attack in step 2 can be represented by some  $\mathcal{U}_1^{\text{acc}}$  on his systems. Alice then sends  $\mathfrak{s}(\mathbf{w}_c)$  and  $\mathfrak{s}(\mathbf{w}_c)$  to Bob, and since Bob passes `TestDummies` in step 5, one can represent his attack by some  $\mathcal{U}_2^{\text{acc}}$  on his systems. Hence, the final state looks like

$$\rho_f^{ABR} = \tilde{\mathcal{E}}' \circ C_{a,s} \circ \mathcal{G} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0 (\rho_0^{ABR}),$$

where  $\tilde{\mathcal{E}}' := \mathcal{U}_2^{\text{acc}} \circ \mathcal{U}_1^{\text{acc}} \circ \tilde{\mathcal{E}}$ .  $\square$



**Lemma G.5.**  $\text{Measure}(\mathfrak{s}(\mathfrak{w}))$ , with  $\mathfrak{w}$  at Alice, is  $\mathcal{G}$ -forcing for any  $\tilde{\mathcal{A}}$  and any  $\tilde{\mathcal{B}}$ , where  $\mathcal{G}$  performs a measurement in the computational basis, and stores the output in registers in both  $A$  and  $B$ .

*Proof.* Since CNOT is forcing, Alice must have actually performed it, and since Bob runs `TestDummies` on  $\mathfrak{s}(\mathfrak{a})$ , Alice must actually send it. Finally, once they measure, the honest party's measurement ensures that the adversary's state also collapses, so that the register only contains the classical measurement result. The same holds for Bob.  $\square$

**Lemma G.6.**  $\text{RGate}(\mathfrak{s}(\mathfrak{w}))$  with  $\mathfrak{w}$  at Alice, is  $\mathcal{G}$ -forcing for any  $\tilde{\mathcal{A}}$  and any  $\tilde{\mathcal{B}}$ , where  $\mathcal{G}(\cdot) := R \cdot R^\dagger$  on  $\mathfrak{w}$ .

*Proof.* This subprotocol is simply a composition of other subprotocols that have already been shown to be forcing.  $\square$

**Lemma G.7.**  $\text{OneQubitClifford}(\mathbb{C}, \mathfrak{s}(\mathfrak{w}))$  with  $\mathfrak{w}$  at Alice, is  $\mathcal{G}$ -forcing for any  $\tilde{\mathcal{A}}$  and any  $\tilde{\mathcal{B}}$ , where  $\mathcal{G}(\cdot) := C \cdot C^\dagger$  on  $\mathfrak{w}$ .  $\square$

*Proof.* Suppose that the state at the beginning in the no-early-abort case is

$$\rho^{AB\mathcal{R}} = \tilde{\mathcal{E}} \circ \mathbb{C}_{b,s} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0 (\rho_0^{AB\mathcal{R}}).$$

Since this subprotocol doesn't actually involve Bob at all, we can simply let  $\tilde{\mathcal{E}}_f$  be the operation applied by Alice on her qubits, and define  $\tilde{\mathcal{E}}' := \tilde{\mathcal{E}}_f \circ \tilde{\mathcal{E}}$ . The final state is therefore

$$\rho_f^{AB\mathcal{R}} = \tilde{\mathcal{E}}' \circ \mathbb{C}_{b,s} \circ \mathcal{G} \circ \mathcal{G}_0 \circ \tilde{\mathcal{E}}_0 (\rho_0^{AB\mathcal{R}}).$$

The same reasoning applies to  $\tilde{\mathcal{B}}$ .  $\square$