

# How to Handle Invalid Queries for Malicious-Private Protocols Based on Homomorphic Encryption

Koji Nuida<sup>12</sup>

<sup>1</sup> Institute of Mathematics for Industry (IMI), Kyushu University, Japan

nuida@imi.kyushu-u.ac.jp

<sup>2</sup> National Institute of Advanced Industrial Science and Technology (AIST), Japan

## Abstract

We consider a setting of two-party computation between a server and a client where every message received by the server is encrypted by a fully homomorphic encryption (FHE) scheme and its decryption key is held by the client only. Akavia and Vald (IACR ePrint Archive, 2021) revisited the privacy problem in such protocols against malicious servers and revealed (as opposed to a naive expectation) that under certain condition, a malicious server can recover the client’s input even if the underlying FHE scheme is IND-CPA secure. They also gave some sufficient conditions for the FHE scheme to ensure the privacy against malicious servers. However, their argument did not consider the possibility that a query from a malicious server to a client involves an invalid ciphertext. In this paper, we show, by giving a concrete example, that if such an invalid query is just rejected by the client, then the sufficient conditions in Akavia and Vald’s result do not in general ensure the privacy against malicious servers. We also propose another option to handle an invalid query in a way that the client returns a random ciphertext (without explicitly rejecting the query), and show that such a protocol is private against malicious servers if the underlying FHE scheme is IND-CPA secure, sanitizable (in the sense of Ducas and Stehlé, EUROCRYPT 2016), and circular secure.

## 1 Introduction

Secure multiparty computation (MPC) is currently one of the most active research areas in cryptography. In MPC, a number of parties have their individual inputs, and by executing a protocol, each party will obtain the own desired output while keeping the own input secret against other parties. (We note that some party may have empty input and may obtain no individual output.) Among various settings and underlying cryptographic tools for MPC, in this paper we focus on secure two-party computation based on homomorphic encryption (HE) in a so-called client-server model. In the model, a party called *client* is supposed to send encrypted input as well as the underlying public key to the other party called *server*, where the public key and the corresponding secret key are generated by the client. Then the server performs some computation using the homomorphic functionality of the encryption scheme and may interact with the client. Finally, the server sends resulting ciphertexts to the client, and the client decrypts them to obtain the output.

In MPC, a *semi-honest* adversary is supposed to completely follow the protocol specification, while a *malicious* adversary may deviate from the protocol specification in an arbitrary manner. In the study of HE-based MPC in the client-server model, usually much attention is paid to the security of a protocol against the (either semi-honest or malicious) client who has the secret key, while less attention seems to be paid to the security against the server. The reason is that in the client-server model every message received by the server is encrypted by the client’s key; in such a case, it is naively expected that the IND-CPA security of the underlying HE scheme ensures at least privacy of the client’s input against the server (though it seems not easy to ensure the correctness of the output against malicious servers). The naive expectation is basically correct when the server has no intermediate interaction with the client. However, protocols in the

client-server model may require intermediate interaction between the server and the client in order to, e.g., realize advanced functionality beyond the power of the underlying additively-HE scheme itself [10, 11, 12] or improve the efficiency of a protocol based on fully homomorphic encryption (FHE) [13]. A recent work by Akavia and Vald [2] revisited this naive expectation for privacy of HE-based protocols against servers with interaction, and revealed that *this expectation may fail* typically for the case of malicious servers and even for some significantly restricted situation with semi-honest servers. They also gave some sufficient conditions for such protocols to have privacy against (semi-honest or malicious) servers, one of which is based on the notion of *sanitization* for ciphertexts of FHE schemes [6] (that is, an operation to make the original FHE scheme circuit-private).

In more detail, the work in [2] focuses on the situation that a server can make a query of the form  $(c, f)$  (or its multi-input and/or multi-output version) to the client where  $c$  is a ciphertext and  $f$  is a function on the plaintext space of the underlying public key encryption (PKE) scheme, and the server then obtains a ciphertext for plaintext  $f(m)$  as the client’s reply where  $m$  is the decryption result of  $c$ . Such a protocol is called an  $(\mathcal{E}, \mathcal{G})$ -aided protocol in [2]. Then in [2], a concrete example of an  $(\mathcal{E}, \mathcal{G})$ -aided protocol was given in which the reply to a kind of query of the form above enables the server to fully recover the client’s input, even though the underlying PKE scheme is IND-CPA secure. Moreover, as a part of the affirmative results in [2], the authors introduced an enhanced variant of IND-CPA security called **funcCPA-security**, and showed that an  $(\mathcal{E}, \mathcal{G})$ -aided protocol ensures privacy against malicious servers if the underlying PKE scheme satisfies **funcCPA-security**. They also showed that **funcCPA-secure** schemes can be constructed based on sanitizable FHE schemes in [6].

However, here we point out that the formulation of **funcCPA-security** as well as the proof of privacy based on **funcCPA-security** in [2] did not consider the possibility that *a query made by the server involves an invalid ciphertext*. As a malicious server may deviate from the protocol specification in an arbitrary manner, such a possibility of invalid queries must be taken into account when the privacy against malicious servers is discussed. The aim of the present paper is to investigate what happens for the results in [2] when invalid queries are considered.

## 1.1 Our Contributions

In this paper, we first consider a seemingly natural option of handling an invalid query (a query involving invalid ciphertext) in an  $(\mathcal{E}, \mathcal{G})$ -aided protocol; the client rejects an invalid query and then even aborts the protocol after such an invalid query (for prohibiting a malicious server from obtaining further information during the protocol). For this option, we show (in Section 4) that an  $(\mathcal{E}, \mathcal{G})$ -aided protocol is *not* always private against malicious servers even if the underlying PKE scheme is **funcCPA-secure** in the sense of [2]. In fact, the “sanitized” version of the FHE scheme FHEW [5] proposed in [6] is **funcCPA-secure** as shown in [2], but our malicious server in Section 4 violates the privacy of an  $(\mathcal{E}, \mathcal{G})$ -aided protocol that rejects invalid queries. Roughly speaking, this is because a certain “shift” of a ciphertext for plaintext 0 in FHEW is still a valid ciphertext (for plaintext 1), while the same shift of a ciphertext for plaintext 1 becomes an invalid ciphertext, therefore those two cases can be distinguished by querying the shifted ciphertext to the client and then observing whether the query is rejected or not.

In more detail, recall that an adversary for **funcCPA-security** defined in [2] is allowed to make some query corresponding to a query in an  $(\mathcal{E}, \mathcal{G})$ -aided protocol. First we extend **funcCPA-security** proposed in [2] in two ways by clarifying how a query involving an invalid ciphertext is handled. In the first option, which we also call **funcCPA-security**, an invalid query is completely prohibited; more precisely, if an invalid query is made by an adversary then the adversary is immediately regarded as losing the security game. In the second option, which we call **funcCPA<sup>†</sup>-security**, if an invalid query is made by an adversary then the query is rejected and the adversary loses the access to the oracle in the sequel. Then we prove (in Theorem 1) that **funcCPA<sup>†</sup>-security** implies **funcCPA-security**, and (in Theorem 3) that any  $(\mathcal{E}, \mathcal{G})$ -aided protocol based on a **funcCPA<sup>†</sup>-secure** encryption scheme is private against malicious servers. On the other hand, **funcCPA-security** in the sense of [2] immediately implies **funcCPA-security** in our sense, but the **funcCPA-security** does not in general imply the privacy of the protocol as the aforementioned counterexample in Section 4 shows.

We note that an essential reason behind the aforementioned counterexample and the non-implication

of privacy from funcCPA-security is that in a PKE scheme, especially for noise-based FHE schemes such as FHEW, it is not always possible to efficiently distinguish invalid ciphertexts from valid ciphertexts by using public information only. In fact, if invalid ciphertexts can be efficiently distinguished by using public information, then funcCPA-security implies funcCPA<sup>†</sup>-security (see Theorem 2) and hence implies privacy of the protocol. Therefore, an important problem is how to ensure the privacy when invalid ciphertexts cannot be efficiently distinguished. To resolve the problem in this case, we introduce another option of handling an invalid query in an  $(\mathcal{E}, \mathcal{G})$ -aided protocol; when the query is invalid, instead of rejecting the query, the client returns a random ciphertext for plaintext 0 to the server (and does not abort the protocol). We call this option *silent abort*. Intuitively, the silent abort is expected to prohibit, owing to the security of the underlying PKE scheme, a malicious server from obtaining information using replies to invalid queries.

Towards proving privacy for  $(\mathcal{E}, \mathcal{G})$ -aided protocols with silent abort based on FHE schemes (as in the proof of privacy in [2] based on FHE), an obstacle is that handing invalid ciphertexts directly is beyond the functionality of FHE schemes (where the behavior of homomorphic evaluation is guaranteed only for valid ciphertexts). To resolve the issue, we utilize the idea of bootstrapping in FHE schemes [8]. Roughly speaking, given a valid or an invalid ciphertext  $c$  as a query, a simulator encrypts each bit of  $c$  to obtain a sequence of ciphertexts  $\tilde{c}_1, \dots, \tilde{c}_\ell$ . Then the simulator applies to  $\tilde{c}_1, \dots, \tilde{c}_\ell$  homomorphic evaluation corresponding to a circuit that distinguishes invalid ciphertexts (by using a secret key), which results in a ciphertext whose plaintext represents whether  $c$  is invalid or not. Such an argument enables us to prove the privacy for  $(\mathcal{E}, \mathcal{G})$ -aided protocols with silent abort based on sanitized FHE schemes, by assuming (as in the circular security) that a public key of the FHE scheme involves encryption of a secret key. More precisely, first we introduce another variant of funcCPA-security corresponding to the case of silent abort, called funcCPA<sup>††</sup>-security, and show (in Theorem 5) that funcCPA<sup>††</sup>-security implies privacy of  $(\mathcal{E}, \mathcal{G})$ -aided protocols with silent abort. Then we prove (in Theorem 6) that a sanitized FHE scheme under the assumption above is funcCPA<sup>††</sup>-secure, hence implying the privacy. It is an open problem to extend the last result above to the case of FHE schemes without circular security, and also to obtain some similar affirmative result in the case of less powerful HE schemes such as additive HE schemes [9] and two-level HE schemes [1, 3, 7].

## 2 Preliminaries

We write  $a \stackrel{\$}{\leftarrow} S$  to mean that an element  $a$  is chosen from a set  $S$  uniformly at random. We let “PPT” stand for “probabilistic polynomial-time”, and consider (unless otherwise noticed) the notion of “polynomial-time” with respect to a security parameter  $\lambda$ . We say that a non-negative function  $f(\lambda)$  in a security parameter  $\lambda$  is *negligible* if for each integer  $k > 0$ , there exists an integer  $N > 0$  satisfying that for any  $\lambda \geq N$  we have  $f(\lambda) \leq 1/\lambda^k$ . We let  $\{0, 1\}^* = \bigcup_{n \geq 0} \{0, 1\}^n$ . For two probability distributions  $X, Y$  over a finite set  $S$ , their statistical distance  $\Delta(X, Y)$  is defined by

$$\Delta(X, Y) = \frac{1}{2} \sum_{s \in S} |\Pr[X = s] - \Pr[Y = s]| .$$

We note that in what follows, any PPT *non-uniform* algorithm can have (in addition to the main input) some auxiliary input called *advice*, which depends solely on the security parameter  $\lambda$  and is sometimes implicit in the notation.

### 2.1 Public Key Encryption

We clarify the notion of public key encryption (PKE) adopted in this paper. In order to concern later whether or not the set of ciphertexts is efficiently recognizable, here we suppose that any ciphertext is encoded as a bit sequence.

**Definition 1** (Public Key Encryption). We say that a tuple of PPT algorithms  $\text{Gen}$ ,  $\text{Enc}$  and a deterministic (i.e., non-probabilistic) polynomial-time algorithm  $\text{Dec}$  is a *public key encryption (PKE) scheme* if it follows the syntax below and satisfies the correctness condition below, where  $\text{PT}$  is the plaintext space (implicitly

depending on the security parameter, the public key, and the secret key),  $\text{CT}_m$  for  $m \in \text{PT}$  are disjoint subsets of  $\{0, 1\}^*$  (implicitly depending on the security parameter, the public key, and the secret key),  $\text{CT} = \bigcup_{m \in \text{PT}} \text{CT}_m$  is the ciphertext space, and  $\text{negl}$  is some negligible function:

- The key generation algorithm  $\text{Gen}$  is given a security parameter  $\lambda$  as input and outputs a pair of a public key  $\text{pk}$  and a secret key  $\text{sk}$ ;  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$ .
- The encryption algorithm  $\text{Enc}$  is given a public key  $\text{pk}$  and a plaintext  $m \in \text{PT}$  as input and outputs an encryption result  $c$ ;  $c \leftarrow \text{Enc}(\text{pk}, m)$ .
- The decryption algorithm  $\text{Dec}$  is given a secret key  $\text{sk}$  and  $c \in \{0, 1\}^*$  as input and outputs a decryption result  $m'$ ;  $m' \leftarrow \text{Dec}(\text{sk}, c)$ . This  $\text{Dec}$  is supposed to satisfy that  $\text{Dec}(\text{sk}, c) = m$  for any  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$ ,  $m \in \text{PT}$ , and  $c \in \text{CT}_m$ .
- (**Correctness for Enc**) With probability at least  $1 - \text{negl}(\lambda)$  for the choice of  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  the following holds: For any  $m \in \text{PT}$ , we have

$$\Pr[\text{Enc}(\text{pk}, m) \in \text{CT}_m] \geq 1 - \text{negl}(\lambda) ,$$

where the latter probability is taken over the internal randomness for  $\text{Enc}$ .

We say that a PKE scheme as above is

- *with privately recognizable ciphertexts* if there exists a deterministic polynomial-time algorithm  $\mathcal{R}$  satisfying that for any  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  and  $c \in \{0, 1\}^*$ , we have  $\mathcal{R}(\text{pk}, \text{sk}, c) = 1$  if  $c \in \text{CT}$  and  $\mathcal{R}(\text{pk}, \text{sk}, c) = 0$  if  $c \notin \text{CT}$ ;
- *with publicly recognizable ciphertexts* if there exists a deterministic polynomial-time algorithm  $\mathcal{R}$  satisfying that for any  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  and  $c \in \{0, 1\}^*$ , we have  $\mathcal{R}(\text{pk}, c) = 1$  if  $c \in \text{CT}$  and  $\mathcal{R}(\text{pk}, c) = 0$  if  $c \notin \text{CT}$ .

For simplifying our argument, we suppose for any PKE scheme that every element of  $\text{PT}$  has the same bit length.

We explain a definition of circular security for PKE schemes, in the sense of IND-WCIRC-CPA<sup>1</sup> security in [4] where we suppose that the secret key is encrypted in a bit-wise manner.

**Definition 2** (Circular Security). We say that a PKE scheme is *circular secure* if it is IND-CPA secure even under the modification that any adversary is also given ciphertexts of every bit of the secret key.

## 2.2 Homomorphic Encryption

We clarify the notion of homomorphic encryption (HE) adopted in this paper, where the plaintext space is fixed to be  $\{0, 1\}$  for simplifying the argument.

**Definition 3** (Homomorphic Encryption). Let  $\mathcal{C} = (\mathcal{C}_\lambda)_{\lambda \geq 1}$  be a family of sets of circuits where, for each circuit in each set  $\mathcal{C}_\lambda$ , its input is a bit sequence and its output is a single bit. We say that  $(\text{Gen}, \text{Enc}, \text{Dec}, \text{Eval})$  is a  *$\mathcal{C}$ -homomorphic encryption ( $\mathcal{C}$ -HE) scheme* (with plaintext space  $\{0, 1\}$ ) if  $(\text{Gen}, \text{Enc}, \text{Dec})$  is a PKE scheme with  $\text{PT} = \{0, 1\}$ ,  $\text{Eval}$  is a PPT algorithm, and it also follows the syntax below and satisfies the correctness and compactness conditions below, where  $\text{negl}$  is some negligible function and  $\rho$  is some positive polynomial:

- Let  $C \in \mathcal{C}_\lambda$  and let  $\ell$  be the input length of the circuit  $C$ . Then the evaluation algorithm  $\text{Eval}$  is given a public key  $\text{pk}$ , the circuit  $C$ , and  $c_1, \dots, c_\ell \in \{0, 1\}^*$  as input and outputs an evaluation result  $c'$ ;  $c' \leftarrow \text{Eval}(\text{pk}, C, c_1, \dots, c_\ell)$ .

- **(Correctness for Eval)** With probability at least  $1 - \text{negl}(\lambda)$  for the choice of  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  the following holds. Let  $C \in \mathcal{C}_\lambda$  and let  $\ell$  be the input length of the circuit  $C$ . Then for any  $m_i \in \{0, 1\}$  ( $i = 1, \dots, \ell$ ) and  $c_i \in \text{CT}_{m_i}$  ( $i = 1, \dots, \ell$ ), we have

$$\Pr[\text{Eval}(\text{pk}, C, c_1, \dots, c_\ell) \in \text{CT}_{C(m_1, \dots, m_\ell)}] \geq 1 - \text{negl}(\lambda) ,$$

where the latter probability is taken over the internal randomness for Eval.

- **(Compactness)** Any ciphertext in CT has bit length at most  $\rho(\lambda)$ .

In the context above, a fully homomorphic encryption (FHE) scheme may mean a  $\mathcal{C}$ -HE scheme with the set  $\mathcal{C}_\lambda$  containing any practically feasible circuits. But here we do not formalize the meaning of “practically feasible” circuits, as our result below is formulated without using the notion of FHE schemes.

### 2.3 Sanitization and Sanitized HE Schemes

The following notion was introduced in [6] (with only slight differences).

**Definition 4** (Sanitization Algorithm). Let  $(\text{Gen}, \text{Enc}, \text{Dec})$  be a PKE scheme. We say that a PPT algorithm Sanitize is a (*statistical*) *sanitization algorithm* if the following conditions hold, where  $\text{negl}$  is some negligible function:

- **(Message-Preserving)** With probability at least  $1 - \text{negl}(\lambda)$  for the choice of  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  the following holds: For any  $m \in \text{PT}$  and  $c \in \text{CT}_m$ , we have

$$\Pr[\text{Sanitize}(\text{pk}, c) \in \text{CT}_m] \geq 1 - \text{negl}(\lambda) .$$

- **(Sanitizing)** With probability at least  $1 - 2^{-\lambda}$  for the choice of  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  the following holds: For any  $m \in \text{PT}$  and  $c_1, c_2 \in \text{CT}_m$ , we have

$$\Delta(\text{Sanitize}(\text{pk}, c_1), \text{Sanitize}(\text{pk}, c_2)) \leq 2^{-\lambda} .$$

Based on the notion of sanitization, a “sanitized” version of HE schemes was introduced in [2]. We explain the construction as follows.

**Definition 5** (Sanitized HE Scheme). Let  $\mathcal{E} = (\text{Gen}, \text{Enc}, \text{Dec}, \text{Eval})$  be a  $\mathcal{C}$ -HE scheme having a sanitization algorithm Sanitize. Then we define the *sanitized scheme*  $\mathcal{E}^{\text{santiz}} = (\text{Gen}, \text{Enc}^{\text{santiz}}, \text{Dec}, \text{Eval}^{\text{santiz}})$  as follows, where Gen and Dec are the same as  $\mathcal{E}$ :

- $\text{Enc}^{\text{santiz}}(\text{pk}, m)$  outputs  $\text{Sanitize}(\text{pk}, \text{Enc}(\text{pk}, m))$ .
- $\text{Eval}^{\text{santiz}}(\text{pk}, C, c_1, \dots, c_\ell)$  computes  $c'_i \leftarrow \text{Sanitize}(\text{pk}, c_i)$  for  $i = 1, \dots, \ell$  and then outputs

$$\text{Sanitize}(\text{pk}, \text{Eval}(\text{pk}, C, c'_1, \dots, c'_\ell)) .$$

It is mentioned in [2] that if  $\mathcal{E}$  is a  $\mathcal{C}$ -HE scheme, then  $\mathcal{E}^{\text{santiz}}$  is also a  $\mathcal{C}$ -HE scheme.

We recall the following example of sanitization algorithms proposed in [6].

**Example 1.** We consider a sanitization algorithm proposed in Section 4.3 of [6], which is applied to an FHE scheme FHEW given in [5]. Let  $q > 0$  be an integer modulus, and let  $n$  be a parameter for dimension. For a secret key  $\vec{s} \in (\mathbb{Z}/q\mathbb{Z})^n$ , a plaintext  $\mu \in \{0, 1\}$ , and some parameters  $\eta$  and  $t$ , we define

$$\text{LWE}_{\vec{s}}^{t:q}(\mu, \eta) = \{(\vec{a}, \langle \vec{a}, \vec{s} \rangle + \mu \cdot \lfloor q/t \rfloor + e) \in (\mathbb{Z}/q\mathbb{Z})^{n+1} : |e| < \eta q\} ,$$

where  $\vec{a} \in (\mathbb{Z}/q\mathbb{Z})^n$  and an element of  $\mathbb{Z}/q\mathbb{Z}$  is identified with an integer in interval  $[-q/2, q/2)$ . The ciphertext space for plaintext  $\mu$  is set to  $\text{CT}_\mu = \text{LWE}_{\vec{s}}^{4:q}(\mu, 1/16)$ ; see [5] and Section 4.3 of [6] for concrete

choices of parameters  $n$  and  $q$ . An outline of the homomorphic operation is as follows. Given two ciphertexts  $c_1 \in \text{CT}_{\mu_1}$  and  $c_2 \in \text{CT}_{\mu_2}$ , the computation part of the operation yields an element of  $\text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4)$  with  $\mu = \text{NAND}(\mu_1, \mu_2)$ , and then the refreshing part converts it back into a ciphertext in  $\text{CT}_{\mu}$ :

$$\text{NAND: } \text{LWE}_{\bar{s}}^{4:q}(\mu_1, 1/16) \times \text{LWE}_{\bar{s}}^{4:q}(\mu_2, 1/16) \rightarrow \text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4) \xrightarrow{\text{Refresh}} \text{LWE}_{\bar{s}}^{4:q}(\mu, 1/16) .$$

On the other hand, a “washing cycle” for elements of  $\text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4)$  was introduced in [6]. An outline is as follows (see Section 4.3 of [6] for the details). It begins with a procedure called homomorphic accumulator to obtain an element of  $\text{LWE}_{\bar{z}}^{2:Q}(\mu, \eta)$  for another key  $\bar{z}$  and certain parameters  $Q$  and  $\eta$ . After this, a rerandomization is performed, by adding to it a random sum of ciphertexts for plaintext 0 and a random noise term. This rerandomization realizes that given two input distributions, the two output distributions have statistical distance that is reduced by a certain factor in comparison to the input distributions. Then the element in  $\text{LWE}_{\bar{z}}^{2:Q}(\mu, \eta)$  is converted back into an element of  $\text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4)$  by a successive application of key switching and modulus switching. This is the process of a single washing cycle:

$$\text{Wash: } \text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4) \rightarrow \text{LWE}_{\bar{z}}^{2:Q}(\mu, \eta) \rightarrow \text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4) .$$

Based on this, the sanitization algorithm for ciphertexts in  $\text{CT}_{\mu}$  is constructed by inserting a sufficient number of washing cycles into the cycle  $\text{LWE}_{\bar{s}}^{4:q}(\mu, 1/16) \rightarrow \text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4) \rightarrow \text{LWE}_{\bar{s}}^{4:q}(\mu, 1/16)$  where the first map is given by  $(\bar{a}, b) \mapsto (2\bar{a}, 2b)$  and the second map is the refreshing procedure used in the construction of NAND:

$$\begin{aligned} \text{Sanitize: } \text{LWE}_{\bar{s}}^{4:q}(\mu, 1/16) &\xrightarrow{\times 2} \text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4) \xrightarrow{\text{Wash}} \text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4) \xrightarrow{\text{Wash}} \dots \\ &\dots \xrightarrow{\text{Wash}} \text{LWE}_{\bar{s}}^{2:q}(\mu, 1/4) \xrightarrow{\text{Refresh}} \text{LWE}_{\bar{s}}^{4:q}(\mu, 1/16) . \end{aligned}$$

## 2.4 Two-Party Protocols

In this paper, we consider two-party protocols  $\pi$  between a client  $\text{Clnt}$  and a server  $\text{Srv}$  where the server has no output, and focus on privacy against the server only. Given input  $x$  for the client  $\text{Clnt}$  and input  $y$  for a possibly malicious server  $\text{Srv}^*$ , let  $\text{view}_{\text{Srv}^*}^{\pi}(x, y)$  denote the random variable of the view (consisting of input, randomness, and received messages during the protocol) for  $\text{Srv}^*$  during the protocol  $\pi$  executed with those inputs. Then the privacy of  $\pi$  against possibly malicious servers is defined as follows:

**Definition 6** (Privacy against Malicious Servers). In the setting above, we say that  $\pi$  is *private against malicious servers* if for any PPT (possibly malicious) server  $\text{Srv}^*$ , there exists a PPT simulator  $\mathcal{S}^*$  satisfying the following: for any PPT non-uniform distinguisher  $\mathcal{D}$ , there exists a negligible function  $\text{negl}$  satisfying that

$$|\Pr[\mathcal{D}(\text{view}_{\text{Srv}^*}^{\pi}(x, y)) = 1] - \Pr[\mathcal{D}(\mathcal{S}^*(1^{\lambda}, y)) = 1]| \leq \text{negl}(\lambda)$$

for any inputs  $x$  and  $y$ .

## 3 Function-Chosen-Plaintext Attack Security, Revisited

In [2], a strengthened variant of CPA security for PKE schemes, called security against function-chosen-plaintext attack (funcCPA-security, in short), was introduced. Let  $\mathcal{E} = (\text{Gen}, \text{Enc}, \text{Dec})$  be a PKE scheme, and let  $\mathcal{G}$  be a family of functions whose input is a number of plaintexts and whose output is a single plaintext. (We note that here we slightly extend the original definition in [2] to allow a multi-input function, though the essence of the definition is not changed.) In order to define the security game, the following decrypt-function-encrypt oracle  $\text{Enc}_{\text{pk}}(\mathcal{G}(\text{Dec}_{\text{sk}}(\cdot)))$  was introduced:

- A query to the oracle  $\text{Enc}_{\text{pk}}(\mathcal{G}(\text{Dec}_{\text{sk}}(\cdot)))$  consists of a function  $g \in \mathcal{G}$  and ciphertexts  $c_1, \dots, c_{\ell} \in \text{CT}$ , where  $\ell$  is the number of input plaintexts for  $g$ . Then the oracle computes  $m_i \leftarrow \text{Dec}(\text{sk}, c_i)$  for  $i = 1, \dots, \ell$ ,  $m' \leftarrow g(m_1, \dots, m_{\ell})$ , and  $c' \leftarrow \text{Enc}(\text{pk}, m')$ , and return  $c'$ .

Roughly speaking,  $\mathcal{E}$  is defined to be **funcCPA**-secure with respect to  $\mathcal{G}$  if  $\mathcal{E}$  is CPA-secure even when the adversary is given access to the oracle above.

Regarding the definition of **funcCPA**-security in [2], here we emphasize that the definition in [2] does not clearly state what happens if an invalid ciphertext is queried to the oracle. By focusing on this point, we introduce the following two variants of **funcCPA**-security with seemingly natural treatment of invalid queries. Intuitively, the first option is that the adversary is simply not allowed to make a query with invalid ciphertext (the adversary loses the security game immediately when such a query is made). The second option is that when a query with invalid ciphertext is made, the oracle aborts and becomes not available in the sequel. Concretely, our security experiment  $\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}}(\lambda)$  with an adversary  $\mathcal{A}$  is defined as follows:

1. The challenger generates a key pair  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  and sets  $\text{abort} = \text{no}$ .
2. The adversary  $\mathcal{A}$  is given  $\text{pk}$  and has access to the oracle  $\text{Enc}_{\text{pk}}(\mathcal{G}(\text{Dec}_{\text{sk}}(\cdot)))$ . Here, when at least one ciphertext queried to the oracle is not in  $\text{CT}$ ,
  - (Case 1) the challenger sets  $\text{abort} = \text{yes}$  and the oracle returns  $\perp$ ;
  - (Case 2) the oracle returns  $\perp$ , and after this,  $\mathcal{A}$  is not allowed to access to the oracle.
3.  $\mathcal{A}$  chooses challenge plaintexts  $m_0, m_1 \in \text{PT}$  and sends them to the challenger.
4. The challenger chooses a challenge bit  $b^* \xleftarrow{\$} \{0, 1\}$ , generates a challenge ciphertext  $c^* \leftarrow \text{Enc}(\text{pk}, m_{b^*})$ , and sends  $c^*$  to  $\mathcal{A}$ .
5.  $\mathcal{A}$  still has access to the oracle  $\text{Enc}_{\text{pk}}(\mathcal{G}(\text{Dec}_{\text{sk}}(\cdot)))$  (unless the oracle access is lost as described above), and finally outputs a bit  $b$ . The output of the experiment is defined to be 1 if  $b = b^*$  and  $\text{abort} = \text{no}$ , and 0 otherwise.

**Definition 7** (**funcCPA**-security). In the setting above, we say that  $\mathcal{E}$  is **funcCPA**-secure (respectively, **funcCPA** $^\dagger$ -secure) with respect to  $\mathcal{G}$  if for any PPT non-uniform adversary  $\mathcal{A}$ , there exists a negligible function  $\text{negl}$  satisfying

$$\Pr[\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}}(\lambda) = 1] \leq \frac{1}{2} + \text{negl}(\lambda) ,$$

where Case 1 (respectively, Case 2) for the behavior with invalid query is adopted.

We note that the **funcCPA**-security and **funcCPA** $^\dagger$ -security become the same as the CPA-security against non-uniform adversaries when  $\mathcal{G} = \emptyset$ .

From now on, we let  $\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda)$  denote the experiment above with Case 2, while  $\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}}(\lambda)$  denotes the experiment with Case 1. We have the following relation for the two versions of security (the converse does not hold in general, as discussed later).

**Theorem 1.** *In the setting above, if  $\mathcal{E}$  is **funcCPA** $^\dagger$ -secure with respect to  $\mathcal{G}$ , then  $\mathcal{E}$  is **funcCPA**-secure with respect to  $\mathcal{G}$ .*

Roughly speaking, this theorem holds because if an adversary wins the security game for **funcCPA**, then an invalid query is not made during the game, in which case the games for **funcCPA** and **funcCPA** $^\dagger$  are equivalent. A formal proof is given below.

*Proof of Theorem 1.* Let  $\mathcal{A}$  be any PPT (non-uniform) adversary for **funcCPA** experiment. We define an adversary  $\mathcal{A}^\dagger$  for **funcCPA** $^\dagger$  experiment as follows:  $\mathcal{A}^\dagger$  internally executes  $\mathcal{A}$ . Here, when  $\mathcal{A}$  generates challenge plaintexts,  $\mathcal{A}^\dagger$  also use them as challenge plaintexts. When  $\mathcal{A}$  produces the final output bit,  $\mathcal{A}^\dagger$  also use it as the output bit. Moreover, when  $\mathcal{A}$  makes an oracle query (in **funcCPA** experiment),  $\mathcal{A}^\dagger$  forwards the query to the oracle (in **funcCPA** $^\dagger$  experiment) and obtains the reply. If this reply is not  $\perp$ , then  $\mathcal{A}^\dagger$  forwards it to  $\mathcal{A}$ . On the other hand, if this reply is  $\perp$ , then  $\mathcal{A}^\dagger$  stops the execution of  $\mathcal{A}$ , sends any two challenge plaintexts to the challenger if  $\mathcal{A}^\dagger$  has not done it, and then outputs a random bit  $b$ .

By definition,  $\mathcal{A}^\dagger$  (as well as  $\mathcal{A}$ ) is PPT and we have

$$\Pr[\text{EXP}_{\mathcal{A}^\dagger, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda) = 1] \geq \Pr[\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}}(\lambda) = 1] .$$

As  $\mathcal{E}$  is  $\text{funcCPA}^\dagger$ -secure, there exists a negligible function  $\text{negl}(\lambda)$  satisfying

$$\Pr[\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}}(\lambda) = 1] \leq \Pr[\text{EXP}_{\mathcal{A}^\dagger, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda) = 1] \leq \frac{1}{2} + \text{negl}(\lambda) .$$

This implies that  $\mathcal{E}$  is  $\text{funcCPA}$ -secure, as desired.  $\square$

The converse of Theorem 1 holds under an additional condition:

**Theorem 2.** *In the setting above, if  $\mathcal{E}$  is  $\text{funcCPA}$ -secure with respect to  $\mathcal{G}$  and is with publicly recognizable ciphertexts, then  $\mathcal{E}$  is  $\text{funcCPA}^\dagger$ -secure with respect to  $\mathcal{G}$ .*

Roughly speaking, this theorem holds because if  $\mathcal{E}$  is with publicly recognizable ciphertexts, then an adversary can avoid making an invalid query to the oracle, in which case the games for  $\text{funcCPA}$  and  $\text{funcCPA}^\dagger$  are equivalent. A formal proof is given below.

*Proof of Theorem 2.* Let  $\mathcal{R}$  be the algorithm implied by the “publicly recognizable ciphertexts” property for  $\mathcal{E}$ . Let  $\mathcal{A}^\dagger$  be any PPT (non-uniform) adversary for  $\text{funcCPA}^\dagger$  experiment. We define an adversary  $\mathcal{A}$  for  $\text{funcCPA}$  experiment as follows:  $\mathcal{A}$  internally executes  $\mathcal{A}^\dagger$ , where the challenge plaintexts and the output bits are common for  $\mathcal{A}$  and  $\mathcal{A}^\dagger$ . Here, when  $\mathcal{A}^\dagger$  makes a query  $(g, c_1, \dots, c_\ell)$  to the oracle (in  $\text{funcCPA}^\dagger$  experiment),  $\mathcal{A}$  first checks if all  $c_i$  are in CT by using the algorithm  $\mathcal{R}(\text{pk}, \cdot)$ . If all  $c_i$  are in CT, then  $\mathcal{A}$  forwards the query to the oracle in  $\text{funcCPA}$  experiment, and sends the reply back to  $\mathcal{A}^\dagger$ . On the other hand, if at least one  $c_i$  is not in CT, then  $\mathcal{A}$  sends  $\perp$  back to  $\mathcal{A}^\dagger$  and in the sequel, rejects any oracle query from  $\mathcal{A}^\dagger$ .

By definition,  $\mathcal{A}$  (as well as  $\mathcal{A}^\dagger$ ) is PPT and we have

$$\Pr[\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}}(\lambda) = 1] = \Pr[\text{EXP}_{\mathcal{A}^\dagger, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda) = 1] .$$

As  $\mathcal{E}$  is  $\text{funcCPA}$ -secure, there exists a negligible function  $\text{negl}(\lambda)$  satisfying

$$\Pr[\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}}(\lambda) = 1] = \Pr[\text{EXP}_{\mathcal{A}^\dagger, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda) = 1] \leq \frac{1}{2} + \text{negl}(\lambda) .$$

This implies that  $\mathcal{E}$  is  $\text{funcCPA}^\dagger$ -secure, as desired.  $\square$

### 3.1 Relation to Two-Party Protocols with Evaluation Queries

Let  $\mathcal{E} = (\text{Gen}, \text{Enc}, \text{Dec})$  and  $\mathcal{G}$  be as above. We consider the class of two-party protocols between a client  $\text{Cnt}$  and a server  $\text{Srv}$  as follows, called  $(\mathcal{E}, \mathcal{G})$ -aided protocols. Here we emphasize that the definition below is almost the same as the one in [2], but in [2], it was not clearly stated how the client should respond when a server’s query involves an invalid ciphertext, which should be clarified as now we are considering malicious servers deviating from the protocol specification. Below we adopt a seemingly natural way of handling such invalid queries, that is to abort the protocol once the client receives an invalid query (to do this efficiently, we assume that the PKE scheme is with privately recognizable ciphertexts in the sense of Definition 1). We also note that the definition in [2] considered only such protocols without server’s input. A protocol in this class is described as follows:

1. **(Client’s input outsourcing phase)** Given client’s input  $\vec{x} = (x_1, \dots, x_{\eta(\lambda)})$ , the client generates  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  and  $e_i \leftarrow \text{Enc}(\text{pk}, x_i)$  for  $i = 1, \dots, \eta(\lambda)$ , and sends  $\text{pk}$  and  $\vec{e} = (e_1, \dots, e_{\eta(\lambda)})$  to the server. (We suppose that the number  $\eta(\lambda)$  of inputs is polynomial-time computable.)

2. **(Server’s computation phase)** The server may interact with the client by sending a query of the form  $(g, \vec{c})$ , where  $g \in \mathcal{G}$  is an  $\ell$ -input function and  $\vec{c} = (c_1, \dots, c_\ell)$ . For such a query,

- if  $c_1, \dots, c_\ell \in \text{CT}$ , then the client computes

$$\text{Enc}(\text{pk}, g(\text{Dec}(\text{sk}, c_1), \dots, \text{Dec}(\text{sk}, c_\ell)))$$

and sends it back to the server;

- otherwise, the client aborts the protocol (the server receives an aborting signal).

Finally, the server sends a message to the client.

3. **(Client’s output phase)** The client generates an output.

Now basically the same argument as Theorem 7 of [2] implies the following:

**Theorem 3.** *Let  $\pi$  be a PPT  $(\mathcal{E}, \mathcal{G})$ -aided protocol as above. If  $\mathcal{E}$  is with privately recognizable ciphertexts and is  $\text{funcCPA}^\dagger$ -secure with respect to  $\mathcal{G}$ , then  $\pi$  is private against malicious servers.*

*Proof.* Let  $\mathcal{R}$  be the algorithm in the definition of “with privately recognizable ciphertexts” property for  $\mathcal{E}$ . Let  $\text{Srv}^*$  be any PPT malicious server. Then we define a simulator  $\mathcal{S}^*$  with server’s input  $y$  as follows:

1.  $\mathcal{S}^*$  generates  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$ , chooses  $x_i^0 \xleftarrow{\$} \text{PT}$  and generates  $e_i^0 \leftarrow \text{Enc}(\text{pk}, x_i^0)$  for  $i = 1, \dots, \eta(\lambda)$ , and sets  $\vec{e}^0 = (e_1^0, \dots, e_{\eta(\lambda)}^0)$ .
2.  $\mathcal{S}^*$  internally executes  $\text{Srv}^*$  with input  $y$  and randomness  $r$  chosen randomly by  $\mathcal{S}^*$ , where the list  $\text{msg}$  is initialized to be empty. Here  $\mathcal{S}^*$  sends the  $\text{pk}$  and  $\vec{e}^0$  to  $\text{Srv}^*$  in Step 1 of the protocol  $\pi$ . When a query  $(g, \vec{c})$  is made by  $\text{Srv}^*$  in Step 2 of  $\pi$ ,  $\mathcal{S}^*$  checks if all  $c_1, \dots, c_\ell$  are in  $\text{CT}$  by using the algorithm  $\mathcal{R}(\text{pk}, \text{sk}, \cdot)$ . If  $c_1, \dots, c_\ell \in \text{CT}$ , then  $\mathcal{S}^*$  computes  $c' \leftarrow \text{Enc}(\text{pk}, g(\text{Dec}(\text{sk}, c_1), \dots, \text{Dec}(\text{sk}, c_\ell)))$ , sends  $c'$  to  $\text{Srv}^*$ , and adds  $c'$  to the list  $\text{msg}$ . On the other hand, if some  $c_i$  is not in  $\text{CT}$ , then  $\mathcal{S}^*$  stops the execution of  $\text{Srv}^*$  and adds the aborting signal to the list  $\text{msg}$ .
3. Finally,  $\mathcal{S}^*$  outputs  $(y, r, \text{msg})$  as the simulated view for  $\text{Srv}^*$ .

We note that  $\mathcal{S}^*$  is a PPT algorithm by the construction.

We assume for the contrary that this simulator does not satisfy the privacy condition. That is, there exists a PPT non-uniform distinguisher  $\mathcal{D}^{(z_\lambda)}$  with advice  $z_\lambda$  satisfying that for any negligible function  $\text{negl}$ , we have

$$\left| \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}, y)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y)) = 1] \right| > \text{negl}(\lambda)$$

for some  $\lambda$  and some input pair  $(\vec{x}, y)$ . In particular, by letting  $(\vec{x}^\lambda, y^\lambda)$  be the pair maximizing the left-hand side above for a fixed  $\lambda$ , it follows that the quantity

$$\left| \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \right| \quad (1)$$

is non-negligible.

Now we define a non-uniform adversary  $\mathcal{A}^+$  against  $\text{funcCPA}^\dagger$ -security for  $\mathcal{E}$  with respect to  $\mathcal{G}$  as follows. The advice for  $\mathcal{A}^+$  is  $(z_\lambda, \vec{x}^\lambda, y^\lambda)$ . Given a public key  $\text{pk}$ ,  $\mathcal{A}^+$  chooses  $I \xleftarrow{\$} \{1, \dots, \eta(\lambda)\}$  and sets the challenge plaintexts to be  $m_0 = x_I^\lambda$  and  $m_1 \xleftarrow{\$} \text{PT}$ . Given a challenge ciphertext  $c^*$ ,  $\mathcal{A}^+$  internally executes  $\text{Srv}^*$  with input  $y^\lambda$  and randomness  $r$  chosen randomly by  $\mathcal{A}^+$ , where  $\text{msg}$  is initialized to be empty. Here  $\mathcal{A}^+$  sends  $\text{pk}$  and  $\vec{e} = (e_1, \dots, e_{\eta(\lambda)})$  to  $\text{Srv}^*$  in Step 1 of  $\pi$ , where  $x_i^0 \xleftarrow{\$} \text{PT}$  and  $e_i \leftarrow \text{Enc}(\text{pk}, x_i^0)$  for  $i = 1, \dots, I - 1$ ,  $e_I = c^*$ , and  $e_i \leftarrow \text{Enc}(\text{pk}, x_i^\lambda)$  for  $i = I + 1, \dots, \eta(\lambda)$ . When a query  $(g, \vec{c})$  is made by  $\text{Srv}^*$  in Step 2 of  $\pi$ ,  $\mathcal{A}^+$  forwards it to the oracle in the  $\text{funcCPA}^\dagger$  game. If the reply is not  $\perp$ , then  $\mathcal{A}^+$  adds the reply to  $\text{msg}$  and sends the reply back to  $\text{Srv}^*$ . If the reply is  $\perp$ , then  $\mathcal{A}^+$  stops the execution of  $\text{Srv}^*$  and adds the aborting

signal to  $\text{msg}$ . Finally,  $\mathcal{A}^+$  executes  $\mathcal{D}^{(z_\lambda)}(y^\lambda, r, \text{msg})$  and outputs the output bit  $b$  of the  $\mathcal{D}^{(z_\lambda)}$ . We note that  $\mathcal{A}^+$  is PPT.

Let  $b^*$  denote the challenge bit in the  $\text{funcCPA}^\dagger$  game. Then the distribution of interaction between  $\mathcal{A}^+$  and  $\text{Srv}^*$  conditioned on index  $I = 1$  and  $b^* = 0$  (respectively,  $I = \eta(\lambda)$  and  $b^* = 1$ ) is identical to the protocol  $\pi$  with malicious server  $\text{Srv}^*$  and input pair  $(\bar{x}^\lambda, y^\lambda)$  (respectively, interaction between the simulator  $\mathcal{S}^*$  and  $\text{Srv}^*$  with input  $y^\lambda$ ). Hence we have

$$\begin{aligned} \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\bar{x}^\lambda, y^\lambda)) = 1] &= \Pr[b = 1 \mid (I, b^*) = (1, 0)] , \\ \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] &= \Pr[b = 1 \mid (I, b^*) = (\eta(\lambda), 1)] . \end{aligned}$$

On the other hand, for each  $2 \leq I_0 \leq \eta(\lambda)$ , the distribution of interaction between  $\mathcal{A}^+$  and  $\text{Srv}^*$  conditioned on  $I = I_0$  and  $b^* = 0$  is identical to that conditioned on  $I = I_0 - 1$  and  $b^* = 1$ . Hence we have

$$\Pr[b = 1 \mid (I, b^*) = (I_0, 0)] = \Pr[b = 1 \mid (I, b^*) = (I_0 - 1, 1)] .$$

Therefore, we have

$$\begin{aligned} &\Pr[b = 1 \mid b^* = 0] - \Pr[b = 1 \mid b^* = 1] \\ &= \sum_{I_0=1}^{\eta(\lambda)} \frac{1}{\eta(\lambda)} \Pr[b = 1 \mid (I, b^*) = (I_0, 0)] - \sum_{I_0=1}^{\eta(\lambda)} \frac{1}{\eta(\lambda)} \Pr[b = 1 \mid (I, b^*) = (I_0, 1)] \\ &= \frac{1}{\eta(\lambda)} (\Pr[b = 1 \mid (I, b^*) = (1, 0)] - \Pr[b = 1 \mid (I, b^*) = (\eta(\lambda), 1)]) \\ &= \frac{1}{\eta(\lambda)} \left( \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\bar{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \right) \end{aligned}$$

and

$$\begin{aligned} &\Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\bar{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \\ &= \eta(\lambda) (\Pr[b = 1 \mid b^* = 0] - \Pr[b = 1 \mid b^* = 1]) \\ &= \eta(\lambda) (1 - \Pr[b = 0 \mid b^* = 0] - \Pr[b = 1 \mid b^* = 1]) \\ &= \eta(\lambda) (1 - 2\Pr[b = 0 \wedge b^* = 0] - 2\Pr[b = 1 \wedge b^* = 1]) \\ &= 2\eta(\lambda) \left( \frac{1}{2} - \Pr[\text{EXP}_{\mathcal{A}^+, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda) = 1] \right) . \end{aligned}$$

As  $\eta(\lambda)$  is polynomially bounded (because the client is PPT) and  $\mathcal{E}$  is  $\text{funcCPA}^\dagger$ -secure, it follows that there exists a negligible function  $\text{negl}_1$  satisfying that the quantity above is at least  $-\text{negl}_1(\lambda)$ .

We also define a PPT adversary  $\mathcal{A}^-$  for  $\mathcal{E}$  in a way that it outputs  $b' = 1 - b$  when  $\mathcal{A}^+$  outputs  $b$ . Then by the argument above, we have

$$\begin{aligned} &\Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\bar{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \\ &= \eta(\lambda) (\Pr[b' = 0 \mid b^* = 0] - \Pr[b' = 0 \mid b^* = 1]) \\ &= \eta(\lambda) (\Pr[b' = 0 \mid b^* = 0] - (1 - \Pr[b' = 1 \mid b^* = 1])) \\ &= \eta(\lambda) (2\Pr[b' = 0 \wedge b^* = 0] + 2\Pr[b' = 1 \wedge b^* = 1] - 1) \\ &= 2\eta(\lambda) \left( \Pr[\text{EXP}_{\mathcal{A}^-, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda) = 1] - \frac{1}{2} \right) . \end{aligned}$$

Again, as  $\mathcal{E}$  is  $\text{funcCPA}^\dagger$ -secure, there exists a negligible function  $\text{negl}_2$  satisfying that the quantity above is at most  $\text{negl}_2(\lambda)$ . Summarizing, we have

$$-\text{negl}_1(\lambda) \leq \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\bar{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \leq \text{negl}_2(\lambda)$$

and the quantity in Eq.(1) is at most the negligible value  $\max\{\text{negl}_1(\lambda), \text{negl}_2(\lambda)\}$ , a contradiction. Hence the simulator  $\mathcal{S}^*$  satisfies the condition showing that the protocol  $\pi$  is private against malicious servers, concluding the proof of Theorem 3.  $\square$

In contrast to Theorem 3, the funcCPA-security of  $\mathcal{E}$  does in general not imply privacy against malicious servers. Intuitively, this is because funcCPA-security completely prohibits an adversary from gaining information by using a query with invalid ciphertexts, while an  $(\mathcal{E}, \mathcal{G})$ -aided protocol as formulated above allows the server to make an invalid query at most once, getting the information that some queried ciphertext is invalid. We provide a concrete example in Section 4.

### 3.2 Relation with Sanitized Schemes

It was shown in Theorem 8 of [2] that if  $\mathcal{E}$  is a CPA-secure  $\mathcal{C}$ -HE scheme with a sanitization algorithm, then the sanitized scheme  $\mathcal{E}^{\text{sanitz}}$  is funcCPA-secure with respect to any function family  $\mathcal{G}$  satisfying that  $\mathcal{G} \subseteq \mathcal{C}_\lambda$  for any  $\lambda$  (where we identify each circuit in  $\mathcal{C}_\lambda$  with the function that determines). Here we note that the definition of funcCPA-security in [2] did not consider the possibility that an invalid ciphertext is queried to the oracle. This point becomes not a problem by choosing an option of not allowing invalid queries, i.e., funcCPA-security (rather than funcCPA $^\dagger$ -security) in our sense. Summarizing, we have the following result in terms of our terminology:

**Theorem 4.** *Let  $\mathcal{E}$  be a CPA-secure  $\mathcal{C}$ -HE scheme with a sanitization algorithm. Let  $\mathcal{G}$  be a function family satisfying that  $\mathcal{G} \subseteq \mathcal{C}_\lambda$  for any  $\lambda$ . Then the sanitized scheme  $\mathcal{E}^{\text{sanitz}}$  is funcCPA-secure with respect to  $\mathcal{G}$ .*

We note that, by combining Theorem 4 with Theorem 2, it follows that if the scheme  $\mathcal{E}$  in Theorem 4 is moreover with publicly recognizable ciphertexts, then  $\mathcal{E}^{\text{sanitz}}$  is funcCPA $^\dagger$ -secure. But in general,  $\mathcal{E}^{\text{sanitz}}$  in Theorem 4 is not funcCPA $^\dagger$ -secure, as shown in Section 4.

## 4 funcCPA-Security Does Not Imply Privacy

In this section, we give a concrete example showing that for an  $(\mathcal{E}, \mathcal{G})$ -aided protocol  $\pi$  as in Section 3.1, the assumption for  $\mathcal{E}$  being (with privately recognizable ciphertexts and) funcCPA-secure with respect to  $\mathcal{G}$  does in general not imply that  $\pi$  is private against malicious servers.

We consider an FHE scheme  $\mathcal{E}$  with sanitization algorithm described in Example 1 of Section 2.3. Its ciphertext space for plaintext  $\mu \in \{0, 1\}$  is  $\text{CT}_\mu = \text{LWE}_{\vec{s}}^{\mathbb{Z}/q\mathbb{Z}}(\mu, 1/16)$  as defined in Section 2.3. This scheme is with privately recognizable ciphertexts; indeed, given an element  $c = (\vec{a}, b) \in (\mathbb{Z}/q\mathbb{Z})^{n+1}$  and a candidate of plaintext  $\mu$ , it is efficient to compute (by using the secret key  $\vec{s}$ ) the term  $e$  satisfying  $b = \langle \vec{a}, \vec{s} \rangle + \mu \cdot \lfloor q/4 \rfloor + e$  and then check if  $|e| < q/16$  or not.

Let the function family  $\mathcal{G}$  be consisting of the identity function  $\text{id}$  only:  $\mathcal{G} = \{\text{id}\}$ . Then by Theorem 4, the sanitized scheme  $\mathcal{E}^{\text{sanitz}}$  is funcCPA-secure with respect to  $\mathcal{G}$ . Now we consider any  $(\mathcal{E}^{\text{sanitz}}, \mathcal{G})$ -aided protocol  $\pi$  in which the server is allowed (by the protocol specification) to make at least one query of the form  $(\text{id}, *)$ . For the protocol  $\pi$  and the underlying scheme  $\mathcal{E}^{\text{sanitz}}$ , we consider a malicious server  $\text{Srv}^*$  that behaves as follows:

- Given a public key and a ciphertext  $c$  of (a part of) the client's input  $x$ ,  $\text{Srv}^*$  computes  $c' = c + (0^n, \lfloor q/4 \rfloor) \in (\mathbb{Z}/q\mathbb{Z})^{n+1}$  where  $0^n = (0, \dots, 0)^T \in (\mathbb{Z}/q\mathbb{Z})^n$  and then makes a query  $(\text{id}, c')$  to the client.

In the situation above, if  $c \in \text{CT}_0$ , i.e.,  $c = (\vec{a}, b)$  and  $b = \langle \vec{a}, \vec{s} \rangle + e$  with  $|e| < q/16$ , then  $c' = (\vec{a}, b')$  with  $b' = \langle \vec{a}, \vec{s} \rangle + \lfloor q/4 \rfloor + e$ , therefore  $c' \in \text{CT}_1$ . In this case, the client on receiving the query above does not abort the protocol. On the other hand, if  $c \in \text{CT}_1$ , i.e.,  $c = (\vec{a}, b)$  and  $b = \langle \vec{a}, \vec{s} \rangle + \lfloor q/4 \rfloor + e$  with  $|e| < q/16$ , then  $c' = (\vec{a}, b')$  with  $b' = \langle \vec{a}, \vec{s} \rangle + 2 \cdot \lfloor q/4 \rfloor + e$ . This  $c'$  is not an element of  $\text{CT} = \text{CT}_0 \cup \text{CT}_1$ , therefore in this case, the client on receiving the query above aborts the protocol. Hence the malicious server can distinguish these two cases, i.e., can know (a part of) the client's input  $x$ , by observing whether the client aborts the

protocol or not. This breaks the privacy of  $\pi$  against malicious servers, despite that  $\mathcal{E}^{\text{sanitz}}$  is (with privately recognizable ciphertexts and) **funcCPA**-secure with respect to  $\mathcal{G}$ .

We also note that in this example, the sanitized scheme  $\mathcal{E}^{\text{sanitz}}$  is not **funcCPA**<sup>†</sup>-secure (as otherwise Theorem 3 would imply that  $\pi$  is private against malicious servers). This also gives a counterexample mentioned at the end of Section 3.2.

## 5 Achieving Privacy against Malicious Servers

We have seen above that in the seemingly natural style of aborting the protocol when receiving an invalid query, the privacy of an  $(\mathcal{E}, \mathcal{G})$ -aided protocol against malicious servers is not always achieved even if the underlying scheme  $\mathcal{E}$  is **funcCPA**-secure and with privately recognizable ciphertexts. In this section, we propose a countermeasure to achieve the privacy against malicious servers with such an underlying scheme  $\mathcal{E}$ . The countermeasure here is to change the client's behavior when receiving an invalid query as explained below. Concretely, we consider a class of protocols described as follows, which we call an  $(\mathcal{E}, \mathcal{G})$ -aided protocol with *silent abort*:

1. **(Client's input outsourcing phase)** Given client's input  $\vec{x} = (x_1, \dots, x_{\eta(\lambda)})$ , the client generates  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$  and  $e_i \leftarrow \text{Enc}(\text{pk}, x_i)$  for  $i = 1, \dots, \eta(\lambda)$ , and sends  $\text{pk}$  and  $\vec{e} = (e_1, \dots, e_{\eta(\lambda)})$  to the server. (We suppose that the number  $\eta(\lambda)$  of inputs is polynomial-time computable.)
2. **(Server's computation phase)** The server may interact with the client by sending a query of the form  $(g, \vec{c})$ , where  $g \in \mathcal{G}$  is an  $\ell$ -input function and  $\vec{c} = (c_1, \dots, c_\ell)$ . For such a query,

- if  $c_1, \dots, c_\ell \in \text{CT}$ , then the client computes

$$\text{Enc}(\text{pk}, g(\text{Dec}(\text{sk}, c_1), \dots, \text{Dec}(\text{sk}, c_\ell)))$$

and sends it back to the server;

- otherwise, the client computes  $\text{Enc}(\text{pk}, 0)$  (where 0 denotes a fixed and polynomial-time computable element of PT) and sends it back to the server.

Finally, the server sends a message to the client.

3. **(Client's output phase)** The client generates an output.

In order to study the privacy of such a protocol, we introduce yet another variant of **funcCPA**-security by also modifying the behavior of the security game when a query is invalid. Let  $\mathcal{E}$  and  $\mathcal{G}$  be as in Section 3. Our security experiment  $\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{fCPA}\ddagger}(\lambda)$  with an adversary  $\mathcal{A}$  is defined as follows:

1. The challenger generates a key pair  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$ .
2. The adversary  $\mathcal{A}$  is given  $\text{pk}$  and has access to the oracle  $\text{Enc}_{\text{pk}}(\mathcal{G}(\text{Dec}_{\text{sk}}(\cdot)))$ . Here, when at least one ciphertext queried to the oracle is not in CT, the oracle returns  $\text{Enc}(\text{pk}, 0)$  to the adversary (where 0 denotes a fixed and polynomial-time computable element of PT).
3.  $\mathcal{A}$  chooses challenge plaintexts  $m_0, m_1 \in \text{PT}$  and sends them to the challenger.
4. The challenger chooses a challenge bit  $b^* \xleftarrow{\$} \{0, 1\}$ , generates a challenge ciphertext  $c^* \leftarrow \text{Enc}(\text{pk}, m_{b^*})$ , and sends  $c^*$  to  $\mathcal{A}$ .
5.  $\mathcal{A}$  still has access to the oracle  $\text{Enc}_{\text{pk}}(\mathcal{G}(\text{Dec}_{\text{sk}}(\cdot)))$ , and finally outputs a bit  $b$ . The output of the experiment is defined to be 1 if  $b = b^*$ , and 0 otherwise.

**Definition 8** (funcCPA<sup>††</sup>-security). In the setting above, we say that  $\mathcal{E}$  is funcCPA<sup>††</sup>-secure with respect to  $\mathcal{G}$  if for any PPT non-uniform adversary  $\mathcal{A}$ , there exists a negligible function  $\text{negl}$  satisfying

$$\Pr[\text{EXP}_{\mathcal{A}, \mathcal{E}, \mathcal{G}}^{\text{funcCPA}^{\dagger\dagger}}(\lambda) = 1] \leq \frac{1}{2} + \text{negl}(\lambda) .$$

Then we give a result analogous to Theorem 3, as follows:

**Theorem 5.** *Let  $\pi$  be a PPT  $(\mathcal{E}, \mathcal{G})$ -aided protocol with silent abort as above. If  $\mathcal{E}$  is with privately recognizable ciphertexts and is funcCPA<sup>††</sup>-secure with respect to  $\mathcal{G}$ , then  $\pi$  is private against malicious servers.*

*Proof.* Let  $\mathcal{R}$  be the algorithm in the definition of “with privately recognizable ciphertexts” property for  $\mathcal{E}$ . Let  $\text{Srv}^*$  be any PPT malicious server. Then we define a simulator  $\mathcal{S}^*$  with server’s input  $y$  as follows:

1.  $\mathcal{S}^*$  generates  $(\text{pk}, \text{sk}) \leftarrow \text{Gen}(1^\lambda)$ , chooses  $x_i^0 \xleftarrow{\$}$  PT and generates  $e_i^0 \leftarrow \text{Enc}(\text{pk}, x_i^0)$  for  $i = 1, \dots, \eta(\lambda)$ , and sets  $\vec{e}^0 = (e_1^0, \dots, e_{\eta(\lambda)}^0)$ .
2.  $\mathcal{S}^*$  internally executes  $\text{Srv}^*$  with input  $y$  and randomness  $r$  chosen randomly by  $\mathcal{S}^*$ , where the list  $\text{msg}$  is initialized to be empty. Here  $\mathcal{S}^*$  sends the  $\text{pk}$  and  $\vec{e}^0$  to  $\text{Srv}^*$  in Step 1 of the protocol  $\pi$ . When a query  $(g, \vec{c})$  is made by  $\text{Srv}^*$  in Step 2 of  $\pi$ ,  $\mathcal{S}^*$  checks if all  $c_1, \dots, c_\ell$  are in CT by using the algorithm  $\mathcal{R}(\text{pk}, \text{sk}, \cdot)$ . If  $c_1, \dots, c_\ell \in \text{CT}$ , then  $\mathcal{S}^*$  computes  $c' \leftarrow \text{Enc}(\text{pk}, g(\text{Dec}(\text{sk}, c_1), \dots, \text{Dec}(\text{sk}, c_\ell)))$ , sends  $c'$  to  $\text{Srv}^*$ , and adds  $c'$  to the list  $\text{msg}$ . On the other hand, if some  $c_i$  is not in CT, then  $\mathcal{S}^*$  computes  $c' \leftarrow \text{Enc}(\text{pk}, 0)$ , sends  $c'$  to  $\text{Srv}^*$ , and adds  $c'$  to the list  $\text{msg}$ .
3. Finally,  $\mathcal{S}^*$  outputs  $(y, r, \text{msg})$  as the simulated view for  $\text{Srv}^*$ .

We note that  $\mathcal{S}^*$  is a PPT algorithm by the construction.

We assume for the contrary that this simulator does not satisfy the privacy condition. That is, there exists a PPT non-uniform distinguisher  $\mathcal{D}^{(z_\lambda)}$  with advice  $z_\lambda$  satisfying that for any negligible function  $\text{negl}$ , we have

$$\left| \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}, y)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y)) = 1] \right| > \text{negl}(\lambda)$$

for some  $\lambda$  and some input pair  $(\vec{x}, y)$ . In particular, by letting  $(\vec{x}^\lambda, y^\lambda)$  be the pair maximizing the left-hand side above for a fixed  $\lambda$ , it follows that the quantity

$$\left| \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \right| \quad (2)$$

is non-negligible.

Now we define a non-uniform adversary  $\mathcal{A}^+$  against funcCPA<sup>††</sup>-security for  $\mathcal{E}$  with respect to  $\mathcal{G}$  as follows. The advice for  $\mathcal{A}^+$  is  $(z_\lambda, \vec{x}^\lambda, y^\lambda)$ . Given a public key  $\text{pk}$ ,  $\mathcal{A}^+$  chooses  $I \xleftarrow{\$} \{1, \dots, \eta(\lambda)\}$  and sets the challenge plaintexts to be  $m_0 = x_I^\lambda$  and  $m_1 \xleftarrow{\$}$  PT. Given a challenge ciphertext  $c^*$ ,  $\mathcal{A}^+$  internally executes  $\text{Srv}^*$  with input  $y^\lambda$  and randomness  $r$  chosen randomly by  $\mathcal{A}^+$ , where  $\text{msg}$  is initialized to be empty. Here  $\mathcal{A}^+$  sends  $\text{pk}$  and  $\vec{e} = (e_1, \dots, e_{\eta(\lambda)})$  to  $\text{Srv}^*$  in Step 1 of  $\pi$ , where  $x_i^0 \xleftarrow{\$}$  PT and  $e_i \leftarrow \text{Enc}(\text{pk}, x_i^0)$  for  $i = 1, \dots, I - 1$ ,  $e_I = c^*$ , and  $e_i \leftarrow \text{Enc}(\text{pk}, x_i^\lambda)$  for  $i = I + 1, \dots, \eta(\lambda)$ . When a query  $(g, \vec{c})$  is made by  $\text{Srv}^*$  in Step 2 of  $\pi$ ,  $\mathcal{A}^+$  forwards it to the oracle in the funcCPA<sup>††</sup> game, adds the oracle’s reply to  $\text{msg}$ , and sends the reply back to  $\text{Srv}^*$ . Finally,  $\mathcal{A}^+$  executes  $\mathcal{D}^{(z_\lambda)}(y^\lambda, r, \text{msg})$  and outputs the output bit  $b$  of the  $\mathcal{D}^{(z_\lambda)}$ . We note that  $\mathcal{A}^+$  is PPT.

Let  $b^*$  denote the challenge bit in the funcCPA<sup>††</sup> game. Then the distribution of interaction between  $\mathcal{A}^+$  and  $\text{Srv}^*$  conditioned on index  $I = 1$  and  $b^* = 0$  (respectively,  $I = \eta(\lambda)$  and  $b^* = 1$ ) is identical to the protocol  $\pi$  with malicious server  $\text{Srv}^*$  and input pair  $(\vec{x}^\lambda, y^\lambda)$  (respectively, interaction between the simulator  $\mathcal{S}^*$  and  $\text{Srv}^*$  with input  $y^\lambda$ ). Hence we have

$$\Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}^\lambda, y^\lambda)) = 1] = \Pr[b = 1 \mid (I, b^*) = (1, 0)] ,$$

$$\Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] = \Pr[b = 1 \mid (I, b^*) = (\eta(\lambda), 1)] .$$

On the other hand, for each  $2 \leq I_0 \leq \eta(\lambda)$ , the distribution of interaction between  $\mathcal{A}^+$  and  $\text{Srv}^*$  conditioned on  $I = I_0$  and  $b^* = 0$  is identical to that conditioned on  $I = I_0 - 1$  and  $b^* = 1$ . Hence we have

$$\Pr[b = 1 \mid (I, b^*) = (I_0, 0)] = \Pr[b = 1 \mid (I, b^*) = (I_0 - 1, 1)] .$$

Therefore, we have

$$\begin{aligned} & \Pr[b = 1 \mid b^* = 0] - \Pr[b = 1 \mid b^* = 1] \\ &= \sum_{I_0=1}^{\eta(\lambda)} \frac{1}{\eta(\lambda)} \Pr[b = 1 \mid (I, b^*) = (I_0, 0)] - \sum_{I_0=1}^{\eta(\lambda)} \frac{1}{\eta(\lambda)} \Pr[b = 1 \mid (I, b^*) = (I_0, 1)] \\ &= \frac{1}{\eta(\lambda)} (\Pr[b = 1 \mid (I, b^*) = (1, 0)] - \Pr[b = 1 \mid (I, b^*) = (\eta(\lambda), 1)]) \\ &= \frac{1}{\eta(\lambda)} \left( \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \right) \end{aligned}$$

and

$$\begin{aligned} & \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \\ &= \eta(\lambda) (\Pr[b = 1 \mid b^* = 0] - \Pr[b = 1 \mid b^* = 1]) \\ &= \eta(\lambda) (1 - \Pr[b = 0 \mid b^* = 0] - \Pr[b = 1 \mid b^* = 1]) \\ &= \eta(\lambda) (1 - 2\Pr[b = 0 \wedge b^* = 0] - 2\Pr[b = 1 \wedge b^* = 1]) \\ &= 2\eta(\lambda) \left( \frac{1}{2} - \Pr[\text{EXP}_{\mathcal{A}^+, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda) = 1] \right) . \end{aligned}$$

As  $\eta(\lambda)$  is polynomially bounded (because the client is PPT) and  $\mathcal{E}$  is  $\text{funcCPA}^\dagger$ -secure, it follows that there exists a negligible function  $\text{negl}_1$  satisfying that the quantity above is at least  $-\text{negl}_1(\lambda)$ .

We also define a PPT adversary  $\mathcal{A}^-$  for  $\mathcal{E}$  in a way that it outputs  $b' = 1 - b$  when  $\mathcal{A}^+$  outputs  $b$ . Then by the argument above, we have

$$\begin{aligned} & \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \\ &= \eta(\lambda) (\Pr[b' = 0 \mid b^* = 0] - \Pr[b' = 0 \mid b^* = 1]) \\ &= \eta(\lambda) (\Pr[b' = 0 \mid b^* = 0] - (1 - \Pr[b' = 1 \mid b^* = 1])) \\ &= \eta(\lambda) (2\Pr[b' = 0 \wedge b^* = 0] + 2\Pr[b' = 1 \wedge b^* = 1] - 1) \\ &= 2\eta(\lambda) \left( \Pr[\text{EXP}_{\mathcal{A}^-, \mathcal{E}, \mathcal{G}}^{\text{fCPA}^\dagger}(\lambda) = 1] - \frac{1}{2} \right) . \end{aligned}$$

Again, as  $\mathcal{E}$  is  $\text{funcCPA}^\dagger$ -secure, there exists a negligible function  $\text{negl}_2$  satisfying that the quantity above is at most  $\text{negl}_2(\lambda)$ . Summarizing, we have

$$-\text{negl}_1(\lambda) \leq \Pr[\mathcal{D}^{(z_\lambda)}(\text{view}_{\text{Srv}^*}^\pi(\vec{x}^\lambda, y^\lambda)) = 1] - \Pr[\mathcal{D}^{(z_\lambda)}(\mathcal{S}^*(1^\lambda, y^\lambda)) = 1] \leq \text{negl}_2(\lambda)$$

and the quantity in Eq.(2) is at most the negligible value  $\max\{\text{negl}_1(\lambda), \text{negl}_2(\lambda)\}$ , a contradiction. Hence the simulator  $\mathcal{S}^*$  satisfies the condition showing that the protocol  $\pi$  is private against malicious servers, concluding the proof of Theorem 5.  $\square$

From now, we aim at extending Theorem 4 to the case of  $\text{funcCPA}^\dagger$ -security (and hence implying the privacy against malicious servers for the protocols), by putting some additional assumptions on the underlying scheme  $\mathcal{E}$ . Let  $\mathcal{E}$  be a  $\mathcal{C}$ -HE scheme with plaintext space  $\text{PT} = \{0, 1\}$ . For simplifying the argument, we suppose that for each fixed security parameter  $\lambda$ , any secret key  $\text{sk}$  has a fixed bit length  $\nu(\lambda)$ , and any ciphertext in  $\text{CT}$  has a fixed bit length  $\rho(\lambda)$ . Then for each function  $g$  in function family  $\mathcal{G}$  that has  $\ell$  inputs from  $\text{PT}$  and has output in  $\text{PT}$ , we consider the function  $\tilde{g}_\lambda$  with  $(\nu(\lambda) + \ell \cdot \rho(\lambda))$ -bit input and one-bit output as follows:

- Given input  $(\tilde{s}, \tilde{c}_1, \dots, \tilde{c}_\ell)$  with  $\tilde{s} \in \{0, 1\}^{\nu(\lambda)}$  and  $\tilde{c}_i \in \{0, 1\}^{\rho(\lambda)}$ , the function  $\tilde{g}_\lambda$  first checks, by regarding  $\tilde{s}$  as a secret key, if each  $\tilde{c}_i$  belongs to CT or not;
  - if  $\tilde{c}_1, \dots, \tilde{c}_\ell \in \text{CT}$ , then  $\tilde{g}_\lambda$  decrypts each  $\tilde{c}_i$  by using  $\tilde{s}$  as a secret key to obtain plaintext  $\tilde{m}_i$ , and finally outputs  $g(\tilde{m}_1, \dots, \tilde{m}_\ell)$ ;
  - otherwise,  $\tilde{g}_\lambda$  outputs 0.

Let  $\tilde{\mathcal{G}}_\lambda = \{\tilde{g}_\lambda \mid g \in \mathcal{G}\}$ . Based on these definitions, we have the following result:

**Theorem 6.** *Let  $\mathcal{E}$  be a CPA-secure  $\mathcal{C}$ -HE scheme as above with a sanitization algorithm. Let  $\mathcal{G}$  be a function family. Suppose that for each  $1 \leq i \leq \nu(\lambda)$ , any public key  $\text{pk}$  for  $\mathcal{E}$  involves an element of  $\text{CT}_{\text{sk}[i]}$ , where  $\text{sk}[i]$  denotes the  $i$ -th bit of the secret key  $\text{sk}$ . (Hence  $\mathcal{E}$  should be circular secure in the sense of Definition 2.) Suppose moreover that  $\tilde{\mathcal{G}}_\lambda \subseteq \mathcal{C}_\lambda$  for any  $\lambda$ . Then the sanitized scheme  $\mathcal{E}^{\text{santz}}$  is  $\text{funcCPA}^{\dagger\dagger}$ -secure with respect to  $\mathcal{G}$ .*

*Proof.* Let  $\tilde{\mathcal{A}}$  be any PPT adversary for the  $\text{funcCPA}^{\dagger\dagger}$ -security game of  $\mathcal{E}^{\text{santz}}$ . Then we define an adversary  $\mathcal{A}$  for the CPA-security game of  $\mathcal{E}$  as follows, where  $\tilde{\text{sk}}_i$  denotes the element of  $\text{CT}_{\text{sk}[i]}$  involved in a public key:

- Given a public key  $\text{pk}$  for  $\mathcal{E}$ ,  $\mathcal{A}$  internally executes  $\tilde{\mathcal{A}}$  with input  $\text{pk}$ .
- When  $\tilde{\mathcal{A}}$  makes a query  $(g, \vec{c})$  with  $\vec{c} = (c_1, \dots, c_\ell)$  in the  $\text{funcCPA}^{\dagger\dagger}$ -security game,
  - if some  $c_i$  is not a  $\rho(\lambda)$ -bit sequence,  $\mathcal{A}$  sends  $\text{Sanitize}(\text{pk}, \text{Enc}(\text{pk}, 0))$  back to  $\tilde{\mathcal{A}}$ ;
  - otherwise,  $\mathcal{A}$  generates  $\tilde{c}_{i,j} \leftarrow \text{Enc}(\text{pk}, c_i[j])$  for each  $1 \leq i \leq \ell$  and  $1 \leq j \leq \rho(\lambda)$  where  $c_i[j]$  denotes the  $j$ -th bit of  $c_i$ ; and  $\mathcal{A}$  computes

$$\text{Sanitize}(\text{pk}, \text{Eval}(\text{pk}, \tilde{g}_\lambda, \tilde{\text{sk}}_1, \dots, \tilde{\text{sk}}_{\nu(\lambda)}, \tilde{c}_{1,1}, \dots, \tilde{c}_{1,\rho(\lambda)}, \dots, \tilde{c}_{\ell,1}, \dots, \tilde{c}_{\ell,\rho(\lambda)}))$$

and sends it back to  $\tilde{\mathcal{A}}$ .

- When  $\tilde{\mathcal{A}}$  generates two challenge plaintexts  $m_0, m_1$ ,  $\mathcal{A}$  sends them to the challenger in the CPA-security game and receives a challenge ciphertext  $c^*$ ; and  $\mathcal{A}$  computes  $\text{Sanitize}(\text{pk}, c^*)$  and sends it back to  $\tilde{\mathcal{A}}$ .
- Finally, when  $\tilde{\mathcal{A}}$  outputs a bit,  $\mathcal{A}$  also outputs the same bit.

Note that  $\mathcal{A}$  is PPT as well as  $\tilde{\mathcal{A}}$ .

For each query  $(g, \vec{c})$  made by  $\tilde{\mathcal{A}}$ , we compare the distribution of the object returned by  $\mathcal{A}$  with the distribution of the reply by the oracle in the true  $\text{funcCPA}^{\dagger\dagger}$ -security game. Let  $E$  denote the event that the key pair  $(\text{pk}, \text{sk})$  chosen in the CPA-security game for  $\mathcal{E}$  does not satisfy one of the ‘‘Sanitizing’’ condition of Definition 4, the ‘‘Correctness’’ condition for  $\text{Enc}$  of Definition 1, and the ‘‘Correctness’’ condition for  $\text{Eval}$  of Definition 3; note that  $E$  occurs with probability at most  $2^{-\lambda} + \text{negl}_1(\lambda) + \text{negl}_2(\lambda)$  where  $\text{negl}_1$  and  $\text{negl}_2$  are negligible functions implied by the correctness of  $\text{Enc}$  and  $\text{Eval}$ , respectively. From now, we suppose that the event  $E$  has not occurred. If some  $c_i$  is not a  $\rho(\lambda)$ -bit sequence, then  $c_i \notin \text{CT}$  and the query is invalid. In this case,  $\tilde{\mathcal{A}}$  in the true  $\text{funcCPA}^{\dagger\dagger}$ -security game receives  $\text{Enc}^{\text{santz}}(\text{pk}, 0) = \text{Sanitize}(\text{pk}, \text{Enc}(\text{pk}, 0))$ , which is identical to the case of the game simulated by  $\mathcal{A}$ .

From now, we focus on the other case where all  $c_1, \dots, c_\ell$  are  $\rho(\lambda)$ -bit sequences. We also focus on the case that  $\tilde{c}_{i,j} \in \text{CT}_{c_i[j]}$  for every  $1 \leq i \leq \ell$  and  $1 \leq j \leq \rho(\lambda)$ , which occurs with probability at least  $1 - \ell \cdot \rho(\lambda) \cdot \text{negl}_1(\lambda)$  due to the correctness of  $\text{Enc}$ . Now suppose that the oracle returns  $\text{Enc}^{\text{santz}}(\text{pk}, \mu)$  in the true  $\text{funcCPA}^{\dagger\dagger}$ -security game. We divide the situation into the two cases below:

- If all  $c_1, \dots, c_\ell$  are elements of CT, then  $\mu = g(\text{Dec}(\text{sk}, c_1), \dots, \text{Dec}(\text{sk}, c_\ell))$ , which is also the same as  $\tilde{g}_\lambda(\text{sk}, c_1, \dots, c_\ell)$  in this case.
- Otherwise,  $\mu = 0$ , which is also the same as  $\tilde{g}_\lambda(\text{sk}, c_1, \dots, c_\ell)$  in this case.

In any case, both  $\text{Enc}(\text{pk}, \mu)$  and

$$\text{Eval}(\text{pk}, \widetilde{g}_\lambda, \widetilde{\text{sk}}_1, \dots, \widetilde{\text{sk}}_{\nu(\lambda)}, \widetilde{c}_{1,1}, \dots, \widetilde{c}_{1,\rho(\lambda)}, \dots, \widetilde{c}_{\ell,1}, \dots, \widetilde{c}_{\ell,\rho(\lambda)})$$

belong to  $\text{CT}_\mu$  with probability at least  $1 - \text{negl}_1(\lambda) - \text{negl}_2(\lambda)$ ; conditioned on this case, the distributions of  $\text{Sanitize}(\text{pk}, \text{Enc}(\text{pk}, \mu)) = \text{Enc}^{\text{santiz}}(\text{pk}, \mu)$  and

$$\text{Sanitize}(\text{pk}, \text{Eval}(\text{pk}, \widetilde{g}_\lambda, \widetilde{\text{sk}}_1, \dots, \widetilde{\text{sk}}_{\nu(\lambda)}, \widetilde{c}_{1,1}, \dots, \widetilde{c}_{1,\rho(\lambda)}, \dots, \widetilde{c}_{\ell,1}, \dots, \widetilde{c}_{\ell,\rho(\lambda)}))$$

have statistical distance at most  $2^{-\lambda}$ .

Summarizing, for each query, the distributions of the objects received by  $\widetilde{\mathcal{A}}$  in the true and the simulated games have statistical distance at most

$$(2^{-\lambda} + \text{negl}_1(\lambda) + \text{negl}_2(\lambda)) + \ell \cdot \rho(\lambda) \cdot \text{negl}_1(\lambda) + (\text{negl}_1(\lambda) + \text{negl}_2(\lambda)) + 2^{-\lambda}$$

which is negligible; we write it as  $\text{negl}_3(\lambda)$ . Moreover, note that  $\mathcal{A}$  wins the CPA-security game if and only if  $\widetilde{\mathcal{A}}$  wins the  $\text{funcCPA}^{\dagger\dagger}$ -security game simulated by  $\mathcal{A}$ . Now by putting the (polynomially bounded) number of queries made by  $\widetilde{\mathcal{A}}$  as  $Q(\lambda)$ , it follows that the output distributions of  $\widetilde{\mathcal{A}}$  in the true and the simulated games have statistical distance at most  $Q(\lambda) \cdot \text{negl}_3(\lambda)$ . Hence we have

$$\Pr[\text{EXP}_{\widetilde{\mathcal{A}}, \mathcal{E}^{\text{santiz}}, \mathcal{G}}^{\text{fCPA}^{\dagger\dagger}}(\lambda) = 1] \leq \Pr[\text{EXP}_{\mathcal{A}, \mathcal{E}}^{\text{CPA}}(\lambda) = 1] + Q(\lambda) \cdot \text{negl}_3(\lambda)$$

where  $\text{EXP}_{\mathcal{A}, \mathcal{E}}^{\text{CPA}}$  denotes the CPA-security experiment for  $\mathcal{E}$  with adversary  $\mathcal{A}$ . As  $\mathcal{E}$  is CPA-secure, we have

$$\Pr[\text{EXP}_{\mathcal{A}, \mathcal{E}}^{\text{CPA}}(\lambda) = 1] \leq \frac{1}{2} + \text{negl}_4(\lambda)$$

for some negligible function  $\text{negl}_4$ . Hence

$$\Pr[\text{EXP}_{\widetilde{\mathcal{A}}, \mathcal{E}^{\text{santiz}}, \mathcal{G}}^{\text{fCPA}^{\dagger\dagger}}(\lambda) = 1] \leq \frac{1}{2} + \text{negl}_4(\lambda) + Q(\lambda) \cdot \text{negl}_3(\lambda)$$

and  $\text{negl}_4(\lambda) + Q(\lambda) \cdot \text{negl}_3(\lambda)$  is negligible. This implies that  $\mathcal{E}^{\text{santiz}}$  is  $\text{funcCPA}^{\dagger\dagger}$ -secure, concluding the proof of Theorem 6.  $\square$

## 6 Conclusion

In this paper, we discussed the privacy problem in FHE-based two-party computation between a malicious server and an honest client, the latter holding the secret key of an underlying FHE scheme. We revisited a previous sufficient condition in [2] to ensure the privacy against the server, and pointed out that their condition is in general not sufficient when the server may make a query to the client with invalid input ciphertexts and the client simply rejects such an invalid query. Then we proposed a modification for their condition assuming that the client does not simply reject any invalid query but responds to such a query by sending a dummy ciphertext back to the server. We also investigated when the condition is indeed satisfied, and proved (roughly speaking) that the condition holds if the underlying FHE scheme is circular secure (in the sense of Definition 2) and sanitizable (in the sense of [6]). It is an open problem to extend the last result above to the case of FHE schemes without circular security, and also to the case of non-fully homomorphic encryption schemes.

**Acknowledgements** This work was supported by JST CREST JPMJCR19F6 and by JSPS KAKENHI Grant Number 19H01109.

## References

- [1] N. Attrapadung, G. Hanaoka, S. Mitsunari, Y. Sakai, K. Shimizu, and T. Teruya, “Efficient Two-level Homomorphic Encryption in Prime-order Bilinear Groups and A Fast Implementation in WebAssembly”, in: Proceedings of ACM ASIACCS 2018, pp.685–697, 2018
- [2] A. Akavia and M. Vald, “On the Privacy of Protocols based on CPA-Secure Homomorphic Encryption”, IACR Cryptology ePrint Archive, Report 2021/803, 2021, Version 20210616:092806 (accessed on October 20, 2021)
- [3] D. Boneh, E.-J. Goh, and K. Nissim, “Evaluating 2-DNF Formulas on Ciphertexts”, in: Proceedings of TCC 2005, LNCS 3378, pp.325–341, 2005
- [4] D. Cash, M. Green, and S. Hohenberger, “New Definitions and Separations for Circular Security”, in: Proceedings of PKC 2012, LNCS 7293, pp.540–557, 2012
- [5] L. Ducas and D. Micciancio, “FHEW: Bootstrapping Homomorphic Encryption in Less Than a Second”, in: Proceedings of EUROCRYPT 2015 (Part I), LNCS 9056, pp.617–640, 2015
- [6] L. Ducas and D. Stehlé, “Sanitization of FHE Ciphertexts”, in: Proceedings of EUROCRYPT 2016 (Part I), LNCS 9665, pp.294–310, 2016
- [7] D. M. Freeman, “Converting Pairing-Based Cryptosystems from Composite-Order Groups to Prime-Order Groups”, in: Proceedings of EUROCRYPT 2010, LNCS 6110, pp.44–61, 2010
- [8] C. Gentry, “Fully Homomorphic Encryption Using Ideal Lattices”, in: Proceedings of ACM STOC 2009, pp.169–178, 2009
- [9] P. Paillier, “Public-Key Cryptosystems Based on Composite Degree Residuosity Classes”, in: Proceedings of EUROCRYPT 1999, pp.223–238, 1999
- [10] K. Shimizu, K. Nuida, and G. Rättsch, “Efficient Privacy-Preserving String Search and an Application in Genomics”, *Bioinformatics*, vol.32, no.11, pp.1652–1661, 2016
- [11] H. Sudo, M. Jimbo, K. Nuida, and K. Shimizu, “Secure Wavelet Matrix: Alphabet-Friendly Privacy-Preserving String Search for Bioinformatics”, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol.16, no.5, pp.1675–1684, 2019
- [12] J. R. Troncoso-Pastoriza, S. Katzenbeisser, and M. Celik, “Privacy Preserving Error Resilient DNA Searching through Oblivious Automata”, in: Proceedings of ACM CCS 2007, pp.519–528, 2007
- [13] W. Wang, Y. Jiang, Q. Shen, W. Huang, H. Chen, S. Wang, X. Wang, H. Tang, K. Chen, K. Lauter, D. Lin, “Toward Scalable Fully Homomorphic Encryption Through Light Trusted Computing Assistance”, arXiv:1905.07766v1, 2019 (accessed on October 20, 2021)