

Security Reductions for White-Box Key-Storage in Mobile Payments

Estuardo Alpirez Bock¹, Chris Brzuska¹, Marc Fischlin², Christian Janson²,
and Wil Michiels^{3,4}

¹ Aalto University

{estuardo.alpirezbock,chris.brzuska}@aalto.fi

² Technische Universität Darmstadt

{marc.fischlin,christian.janson}@cryptoplexity.de

³ Technische Universiteit Eindhoven

⁴ NXP Semiconductors

wil.michiels@nxp.com

Abstract. The goal of white-box cryptography is to provide security even when the cryptographic implementation is executed in adversarially controlled environments. White-box implementations nowadays appear in commercial products such as mobile payment applications, e.g., those certified by Mastercard. Interestingly, there, white-box cryptography is championed as a tool for secure storage of payment tokens, and importantly, the white-boxed storage functionality is bound to a hardware functionality to prevent code-lifting attacks.

In this paper, we show that the approach of using hardware binding and obfuscation for secure storage is conceptually sound. Following security specifications by Mastercard, we first define security for a *white-box key derivation functions* (WKDF) that is bound to a hardware functionality. WKDFs with hardware-binding model a secure storage functionality, as the WKDFs in turn can be used to derive encryption keys for secure storage. We then provide a proof-of-concept construction of WKDFs based on pseudorandom functions (PRF) and obfuscation. To show that our use of cryptographic primitives is sound, we perform a cryptographic analysis and reduce the security of our WKDF to the cryptographic assumptions of indistinguishability obfuscation and PRF-security. The hardware-functionality that our WKDF is bound to is a PRF-like functionality. Obfuscation helps us to hide the secret key used for the verification, essentially emulating a signature functionality as is provided by the Android key store.

We rigorously define the required security properties of a *hardware-bound white-box payment application* (WPAY) for generating and encrypting valid payment requests. We construct a WPAY, which uses a WKDF as a secure building block. We thereby show that a WKDF can be securely combined with *any* secure symmetric encryption scheme, including those based on standard ciphers such as AES.

1 Introduction

Near-field communication (NFC) protocols have opened up new possibilities for mobile payment applications, such as those offered by Mastercard, Visa, or Google wallet [Sma14]. Traditionally, the NFC traffic was processed by a secure hardware component in mobile devices, that performed cryptographic operations. In 2015, Android 4.4 introduced Host Card Emulation (HCE), which allows the application processor of a mobile device to use the NFC communication, too. In this case, the cryptographic functions of mobile applications are implemented *software-only*, increasing flexibility and device coverage of the applications since no secure hardware element is required. However, implementing cryptographic functions in software leads to new attack vectors.

White-box cryptography. One of the core cryptographic protection technologies for HCE (as listed by the Smart Card Alliance Mobile & NFC Council [Sma14]) is the use of *white-box cryptography*. Cryptography in the white-box attack model was introduced by Chow, Eisen, Johnson, and van Oorschot in 2002 (CEJO [CEJvO03,CEJv03]). In the white-box attack scenario, an adversary has access to the program code and the complete execution environment of a cryptographic implementation, and the goal of white-box cryptography is to remain secure despite such strong attack capabilities. Unlike in the Digital Rights Management scenario considered by CEJO where the user is considered as an adversary, in an HCE context, the goal of white-box cryptography is protect an honest user against attacks performed on their device, e.g., by malware that can observe program code and executions.

Commercial payment applications. Payment applications need to store secret information such as transaction tokens that are decrypted when a transaction is performed. In the absence of a secure element, the tokens are stored in insecure memory and likewise, the decryption operations are performed by an insecure CPU. Thus, to protect against adversaries that use their access to the storage/CPU to extract secret information and perform payment transactions on their own, over the past years, white-box cryptography has been broadly adopted by those offering commercial payment applications. The Mastercard security guidelines for payment applications, for instance, make the use of white-box cryptography *mandatory* for implementing storage protection in order to achieve an advanced security level (see *Local Database Encryption*, Chapter 5 in [Mas17]).

For a successful payment, the user's device needs to be close to an NFC reader. An attacker on a user's device can thus only alter a payment a user aims to make, but cannot make payments independently at readers of their choice,⁵ unless the attacker gains independence from the user's device. The attacker could gain independence by (a) extracting the key, or (b) performing a code-lifting

⁵ We discuss relay attacks later.

attack [Wys12]. Thus, white-box cryptography in commercial payment applications needs to achieve *hardware-binding* (also see [CdRP14,SdHM15,BBIJ17] for discussions of usefulness). As a consequence, commercial applications implement white-box cryptography with a hardware anchor, essentially reaching a middle-ground between software-only and hardware-only security for cryptographic implementations.

Hardware-binding on Android. The Mastercard guidelines (see Chapter 5 in [Mas17]) recommend, to the very least, to use a unique device fingerprint for device identification. For instance, Android allows to perform checks on identifiers such as the hardware serial number, the ESN (electronic serial number) or IMEI (international mobile equipment identity) of the device via its `Build` and `TelephonyManager` classes [And18a,And18c]. This technique helps to mitigate code-lifting attacks as long as the value remains secret and/or interception of this value between hardware and software can actually be prevented. For an advanced security level, however, the guidelines suggest the use of the functionalities of the Android Key Store. The Android Key Store, e.g., implements RSA signatures, and relies on whatever secure hardware features the Android device provides. Signatures are a more useful binding functionality than single identifier values, since for each new input, they provide a different output.

Conceptual validation. In this paper, we show that the wide-spread practical approach for building secure payment applications based on white-box cryptography is conceptually sound. We split our study into two parts:

- a hardware-bound white-box key derivation function (WKDF) which provides (1) hardware-binding and (2) secure storage;
- a secure payment application that performs (3) symmetric encryption of data on top of the WKDF.

Note that the Mastercard guidelines merely specify best practices but omit a design blueprint. Our goal is to explicate how exactly a sound design shall proceed, and what security properties the underlying primitives should obey.

Hardware-bound white-box key derivation function (WKDF). Our WKDF notion builds on top of a standard (black-box) key derivation function (KDF). We here consider a lightweight notion for the KDF that takes uniformly random keys and a second, non-random value as input, and returns pseudorandom keys of fixed length. We therefore use the terms KDF and pseudorandom function (PRF) interchangeably, abstracting away additional KDF features such as varying output lengths (c.f. Krawczyk [Kra10]). We introduce the IND-WKDF security notion for WKDFs that models the previously discussed white-box attack scenario. I.e., the adversary is given full access to the white-box implementation of the WKDF as well as limited access to the hardware. If the adversary uses its hardware access, the adversary is able to evaluate the WKDF, but if

the adversary has no access to the relevant hardware values, e.g., carrying out a code-lifting attack, then the adversary learns nothing about the WKDF values, which is modeled by a real-or-random oracle for derived keys.

Existing security notions for white-box cryptography. Outside of payment applications, Delerablée, Lepoint, Paillier, and Rivain (DLPR [DLPR14]) defined the meanwhile popular notion of *incompressibility* (for constructions and definitional refinements see [DLPR14,BI15,BIT16,BBK14,FKKM16,BAB⁺19]). Here, an adversary shall not be able to compress a cryptographic program without losing part of its functionality. Incompressibility seems unrelated to achieving the goal of real-or-random key indistinguishability for derived keys, as we aim for the WKDF in our application. Additionally, DLPR define *traceability* which is a security notion intended to trace malicious users (typically in a DRM setting) if they illegally share their white-box program with others. Traceability is a helpful to mitigate code-lifting and re-distribution attacks in the case of malicious users, but not helpful to protect honest users from adversaries copying and misusing their software.

Finally, DPLR also discuss one-wayness and security against key extraction, a baseline security property for white-box cryptography, which are both implied by our IND-WKDF notion for WKDF. Namely, if an adversary can extract the key, then the adversary can evaluate the KDF on all points itself and thereby distinguish derived keys from random keys. Similarly, an adversary can use an inversion algorithm to distinguish real derived keys from random values. Thus, a IND-WKDF-secure WKDF also resists key extraction and inversion attacks.

Secure payment application. We introduce a secure hardware-bound payment application scheme (WPAY). Its basic functionality is to encrypt and to authenticate valid payment requests to a server. We model validity by a predicate that acts as filter function. E.g., the filter could only allow for certain date ranges or limits the upper bound on the payment, while the server generically accepts payments of arbitrary amounts and ranges.

Our security notion IND-WPAY gives the adversary the white-box payment application WPAY. The adversary can query a hardware oracle that provides them with the necessary hardware values to generate a request using WPAY. This models that the adversary can observe (and interfere with) honest user evaluations. As soon as the adversary loses access to the hardware, confidentiality and integrity of the user requests should hold. IND-WPAY models both properties via an indistinguishability game. Note that in particular, IND-WPAY captures code-lifting attacks. Namely, the adversary has access to WPAY throughout the experiment, but only limited access to the hardware. IND-WPAY models that in the absence of the hardware, no valid requests can be generated, even given WPAY.

Constructions. To instantiate our approach for building a WKDF, we first need to specify a hardware functionality. One idea could be to rely on a signa-

ture functionality as provided for example by the Android Key Store. I.e., WKDF would send a request to the hardware, the hardware signs it, and WKDF then verifies the signature with the public verification key. But we need to (1) hide the software-related key of our WKDF and (2) make it inseparable from the verification algorithm that checks hardware values, which is both achieved by applying indistinguishability obfuscation techniques. This, in turn, forces us to use puncturable primitives for the security reduction to work. One option could thus be to use the puncturable signature scheme by Bellare, Stepanovs and Waters [BSW16] which, notably, itself is based on indistinguishability obfuscation. To avoid this double form of obfuscation for the construction, one layer for the puncturable signature scheme and one for the hiding and binding of our KDF key, we instead use a faster symmetric-key primitive in form of a (puncturable) PRF (essentially as a message authentication code). This puncturable PRF is obfuscated once within the hardware-linked KDF construction to ensure the required security.

Hence, we build a WKDF and prove its IND-WKDF security, following techniques by Sahai and Waters [SW14]. This construction assumes puncturable PRFs (which are equivalent to one-way functions) and indistinguishability obfuscation. Given an IND-WKDF-secure WKDF we then prove that another layer of indistinguishability obfuscation can be used to bind the WKDF to an arbitrary secure symmetric encryption scheme and a filter function to obtain an IND-WPAY-secure white-box payment application WPAY.

Discussion and limitations. Note that our constructions are conceptual validations and not practically efficient due to the tremendous inefficiency of indistinguishability obfuscation (see [AHKM14]). In practice, the obfuscation need to be implemented by a mix of efficient obfuscation techniques, combined with practical white-box techniques, e.g. [GPRW18]. Thus, our work does not allow to immediately bypass the difficulty of building white-box implementations —as apparent in the past white-box competitions [ECR17,cyb19], where only three design candidates submitted towards the end of the second competition remained unbroken— in practice. However, our theoretical feasibility result allows us to conclude that secure white-box implementations based on strong cryptographic assumptions —indistinguishability obfuscation is not yet a mature cryptographic primitive— are indeed possible. Our results not only affirm that building secure white-box cryptography is possible, but they also explain *how* such a secure white-box implementation can be designed.

As efficiency of indistinguishability obfuscation has not yet reached reasonable levels, let us now discuss security and efficiency of current practical white-box implementations. E.g., the winner of the first white-box competition [ECR17] had a binary size of 17MB and needed 0.37 seconds for an encryption which is reasonably close to practical needs. It was broken eventually, but resisted key extraction attacks for up to 28 days [Riv17]. This temporary robustness turns out to be indeed useful. Namely, in practice, the goal is to maintain a complexity gap between the effort of the attacker and the effort of

the designer. As software can be updated in a regular interval, one can achieve a reasonable practical security level by replace a white-box implementation, in each update, by a newer generation.

Two considerations are important to be taken into account: (1) The white-box implementation should not be susceptible to (variants of) automated attacks, since these can be implemented with little effort, see [ABBB⁺19]. (2) Reverse-engineering efforts against previous generations of white-box implementations shall not help the attacker against the new generation of the white-box implementation. I.e., the designer needs to come up with a paradigm that allows to systematically inject a certain amount of creativity into the system that needs to be reverse-engineered anew each time. For example, if security requires bootstrapping security from white-boxing another cryptographic primitive such as a PRF, then one can use a different PRF each time, harvesting the large cryptographic research of PRF constructions. In our model, we bind each payment application to a fresh hardware sub key, derived from a master hardware key, and we allow the adversary to see other derived hardware keys. This models that potentially, earlier construction might have been broken and might have revealed the derived hardware key in use. An alternative is to directly bind to a signature functionality so that revealing the verification algorithm does not constitute an attack vector.

Note that our models do not consider plain relay attacks [Han05,FHMM11] where the adversary forwards the intended communication without altering it. These attacks need to be prevented by other means, e.g., via distance-bounding protocols [ABB⁺19], or at least mitigated via heuristics such as location correlation between phone and NFC reader [BKTS16]. Note however that we of course capture attacks where, say, the adversary modifies user requests or tries to create new requests by himself.

2 Preliminaries and Notation

By $a \leftarrow A(x)$, we denote the execution of a deterministic algorithm A on input x and the assignment of the output to a , while $a \leftarrow_s A(x)$ denotes the execution of a randomized algorithm and the assignment of the output to a . We denote by $:=$ the process of initializing a set, e.g. $X := \emptyset$. By $x \leftarrow_s X$ we denote the process of randomly and uniformly sampling an element x from a given set X . Slightly abusing notation, we also use $x \leftarrow_s X$ to denote the sampling of x according to probability distribution X . We then denote the probability that the event $E(x)$ happens by $\Pr_{x \leftarrow_s X}[E(x)]$ or sometimes simply $\Pr[E(x)]$. We write oracles as superscript to the adversary \mathcal{A}^O . In cases when an adversary is granted access to a larger number of oracles, we write oracles also as subscript to the adversary $\mathcal{A}_{O_3, O_4}^{O_1, O_2}$. PPT denotes probabilistic polynomial-time and $\text{poly}(n)$ is an unspecified polynomial in the security parameter. Note that all algorithms receive the security parameter 1^n in unary notation as input implicitly. We write it explicitly only occasionally for clarity. We now review useful definitions, starting with nonce-based encryption, see Rogaway [Rog04].

Definition 1 (Symmetric Encryption). A nonce-based symmetric encryption scheme SE consists of a pair of deterministic polynomial-time algorithms (Enc, Dec) with the syntax $c \leftarrow \text{Enc}(k, m, nc)$ and $m/\perp \leftarrow \text{Dec}(k, c, nc)$. The algorithm Enc takes as input a randomly generated key k of length n , a nonce nc , a message m , and outputs a ciphertext c . Dec takes as input a randomly generated key k of length n , a nonce nc , a ciphertext c , and outputs either a message m or an error symbol \perp . Moreover, the encryption scheme SE satisfies correctness, if for all nonces $nc \in \{0, 1\}^n$ and for all messages $m \in \{0, 1\}^*$,

$$\Pr[\text{Dec}(k_{\text{SE}}, \text{Enc}(k_{\text{SE}}, m, nc), nc) = m] = 1$$

where the probability is over sampling k .

Definition 2 and Figure 1 specify the security of an *authenticated encryption scheme* [BN00,Rog02]. Here, the adversary is provided with a left-or-right encryption oracle and a decryption oracle where it can submit arbitrary ciphertexts except for challenge ciphertexts obtained from the encryption oracle. If $b = 0$, the decryption oracle is functional. If $b = 1$, the decryption oracle always returns \perp which models ciphertext integrity. In the security game, we use **assert** as a shorthand to say that if the **assert** condition is violated, then the oracle returns an error symbol \perp .

Definition 2 (AE-security). A nonce-based symmetric encryption scheme $\text{SE} = (\text{Enc}, \text{Dec})$ is called an authenticated encryption scheme or *AE-secure* if all PPT adversaries \mathcal{A} have negligible distinguishing advantage in the game $\text{Exp}_{\text{SE}, \mathcal{A}}^{\text{AE}}(1^n)$, specified in Figure 1.

Note that we demand the authenticated encryption scheme to be deterministic because we will later execute the algorithm in an untrusted environment and cannot count on strong randomness. This, in turn, implies that we cannot allow the adversary to re-use any of the previous queries (m, nc) , or else it would be easy to determine b from two queries (m_0, m_1, nc) and (m_0, m'_1, nc) .

We provide a formal definition of a key derivation function that produces pseudorandom keys. Note that our definition corresponds to a PRF, i.e., it is highly simplified compared to the framework of Krawczyk [Kra10]. In our definition a key derivation function takes as input a key k_{kdf} , a context string e

$\text{Exp}_{\text{SE}, \mathcal{A}}^{\text{AE}}(1^n)$	$\text{OENC}(m_0, m_1)$	$\text{ODEC}(nc, c)$
$b \leftarrow_{\$} \{0, 1\}$	assert $ m_0 = m_1 $	assert $c \notin C$
$k_{\text{SE}} \leftarrow_{\$} \{0, 1\}^n$	$nc \leftarrow_{\$} \{0, 1\}^n$	if $b = 1$ then
$b' \leftarrow_{\$} \mathcal{A}^{\text{OENC}, \text{ODEC}}(1^n)$	$c \leftarrow \text{Enc}(k_{\text{SE}}, m_b, nc)$	return \perp
return $(b' = b)$	$C := C \cup \{c\}$	else
	return (nc, c)	return $m \leftarrow \text{Dec}(k_{\text{SE}}, c, nc)$

Fig. 1. The $\text{Exp}_{\text{SE}, \mathcal{A}}^{\text{AE}}(1^n)$ security game.

as well as the security parameter 1^n . In comparison to Krawczyk's definition, we simplify the presentation and omit the details of the smoothing step turning raw key material into random strings, the salting, and the length parameter, assuming that the key k_{kdf} is already appropriate and the length of the returned key is equal to $|k_{\text{kdf}}|$.

$\text{Exp}_{\text{KDF}, \mathcal{A}}^{\text{IND-KDF}}(1^n)$	$\text{OKDF}(e)$
$b \leftarrow_{\mathcal{S}} \{0, 1\}$	if $e \notin \mathcal{Q}$
$\mathcal{Q} := \emptyset$	$\mathcal{Q} := \mathcal{Q} \cup \{e\}$
$k_{\text{kdf}} \leftarrow_{\mathcal{S}} \text{Kgen}(1^n)$	if $b = 1$
$b' \leftarrow_{\mathcal{S}} \mathcal{A}^{\text{OKDF}}(1^n)$	$\hat{k} \leftarrow \text{KDF}(k_{\text{kdf}}, e)$
if $b' = b$	else
return 1	$\hat{k} \leftarrow_{\mathcal{S}} \{0, 1\}^n$
else	return \hat{k}
return 0	else
	return \perp

Fig. 2. $\text{Exp}_{\text{KDF}, \mathcal{A}}^{\text{IND-KDF}}(1^n)$ security game

$\text{Exp}_{\text{PPRF}, \mathcal{A}}^{\text{IND-PPRF}}(1^n)$
$b \leftarrow_{\mathcal{S}} \{0, 1\}$
$k \leftarrow_{\mathcal{S}} \{0, 1\}^n$
$(z, \text{state}) \leftarrow_{\mathcal{S}} \mathcal{A}(1^n)$
$k_z \leftarrow_{\mathcal{S}} \text{Punct}(k, z)$
if $b = 1$ then
$y \leftarrow \text{PPRF}(k, z)$
else $y \leftarrow_{\mathcal{S}} \{0, 1\}^n$
$b' \leftarrow_{\mathcal{S}} \mathcal{A}(k_z, y, \text{state})$
return $(b' = b)$

Fig. 3. IND-PPRF security game

Definition 3 (Key Derivation Function). A KDF scheme consists of a randomized key generation algorithm Kgen and a key derivation function KDF that is a deterministic algorithm that takes as input a key $k_{\text{kdf}} \leftarrow_{\mathcal{S}} \text{Kgen}(1^n)$ and a context string e . The KDF returns a key \hat{k} of length $|k_{\text{kdf}}|$.

Definition 4 (IND-KDF-security). A key derivation function KDF is said to be IND-KDF-secure if all PPT adversaries \mathcal{A} have negligible distinguishing advantage in game $\text{Exp}_{\text{PPRF}, \mathcal{A}}^{\text{IND-PPRF}}(1^n)$, specified in Figure 2.

Next we present the definition of a length-doubling pseudorandom generator.

Definition 5 (Pseudorandom Generator). A deterministic, polynomial-time computable function $\text{PRG} : \{0, 1\}^* \rightarrow \{0, 1\}^*$ is a pseudorandom generator if:

- **Length-expansion:** For all $x \in \{0, 1\}^*$, $|\text{PRG}(x)| = 2|x|$.
- **Pseudorandomness:** For all PPT \mathcal{A} , $\text{Adv}_{\text{PRG}, \mathcal{A}}(n) :=$

$$|\Pr_{x \leftarrow_{\mathcal{S}} \{0, 1\}^n} [\mathcal{A}(\text{PRG}(x)) = 1] - \Pr_{z \leftarrow_{\mathcal{S}} \{0, 1\}^{2n}} [\mathcal{A}(z) = 1]|$$

is negligible in n .

We define pseudorandom functions with identical input, output and key length.

Definition 6 (Pseudorandom Function). A deterministic, polynomial-time computable function $\text{PRF} : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^n$ for all $n \in \mathbb{N}$, is a pseudorandom function if for all PPT \mathcal{A} , $\text{Adv}_{\mathcal{A}, \text{PRF}}(n) :=$

$$\left| \Pr_{k \leftarrow \mathbb{S}} \{0, 1\}^n \left[\mathcal{A}^{\text{PRF}(k, \cdot)}(1^n) = 1 \right] - \Pr_{F \leftarrow \mathbb{S}} \{G : \{0, 1\}^n \rightarrow \{0, 1\}^n\} \left[\mathcal{A}^{F(\cdot)}(1^n) = 1 \right] \right|$$

is negligible in n .

Puncturable PRFs (PPRF) were introduced by Boneh and Waters [BW13]. PPRFs have a punctured key which allows to evaluate the PPRF on all inputs, except for one where the function still looks random.

Definition 7 (PPRF). A puncturable pseudorandom function scheme PPRF consists of a triple $(\text{PPRF}, \text{Punct}, \text{Eval})$, which are defined as follows:

- $\text{PPRF}(k, x)$: This is a standard PRF evaluation algorithm. As before, this deterministic polynomial-time algorithm takes as input a key k and input x , both of length n and returns a value y of length n .
- $\text{Punct}(k, z)$: This PPT algorithm takes as input a key $k \in \{0, 1\}^n$ and an input value $z \in \{0, 1\}^n$. It outputs a punctured key $k_z \leftarrow \mathbb{S} \text{Punct}(k, z)$.
- $\text{Eval}(k_z, x)$: This deterministic polynomial-time algorithm takes as input a punctured key k_z and some input $x \in \{0, 1\}^n$ and returns \perp if $x = z$, and a value $y \in \{0, 1\}^n$ otherwise.

A puncturable PRF is said to be correct, if for all security parameter n , all $k \in \{0, 1\}^n$, every value $z \in \{0, 1\}^n$ and all $x \in \{0, 1\}^n$, $x \neq z$, it holds that

$$\Pr[\text{Eval}(\text{Punct}(k, z), x) = \text{PPRF}(k, x)] = 1.$$

PPRF security requires that the PPRF value on k and z is indistinguishable from random, even when given the punctured key k_z . Note that for our purposes, we only rely on security for *random* inputs rather than adversarially chosen ones, i.e., we use a less powerful assumption which makes our result stronger.

Definition 8 (IND-PPRF-security). A PPRF scheme is said to be IND-PPRF-secure if all probabilistic polynomial-time adversaries \mathcal{A} have negligible distinguishing advantage in the IND-PPRF game defined in Figure 3.

An indistinguishability obfuscator (iO) ensures that the obfuscation of any two functionally equivalent programs (i.e. circuits) are computationally indistinguishable. In the following definition, a distinguisher \mathcal{D} is an adversary that aims at identifying which of the two programs has been obfuscated.

Definition 9 (Admissible Circuit Sampler). Let p be a polynomial. A PPT algorithm \mathcal{S} is called a p -admissible circuit sampler if

$$1 - \Pr_{(C_0, C_1) \leftarrow \mathbb{S}} \mathcal{S}(1^n) [\forall x \in \{0, 1\}^n C_0(x) = C_1(x)]$$

is negligible in n and for all $n \in \mathbb{N}$ and all pairs (C_0, C_1) in the range of $\mathcal{A}(1^n)$, it holds that the size of C_0 and the size of C_1 is upper bounded by $p(n)$.

Definition 10 (Indistinguishability Obfuscator). A PPT algorithm iO , parametrized by a polynomial p , is called an indistinguishability obfuscator if for any p -admissible circuit sampler \mathcal{S} the following conditions are satisfied:

Correctness. For all circuits C and for all inputs x to the circuit,

$$\Pr[C'(x) = C(x) : C' \leftarrow_{\mathcal{S}} \mathsf{iO}(C)] = 1,$$

where the probability is over the randomness of iO .

Security. For all p -admissible \mathcal{S} and any PPT distinguisher \mathcal{D} , the following distinguishing advantage is negligible:

$$\left| \Pr_{(C_0, C_1) \leftarrow_{\mathcal{S}}} [\mathcal{D}(\mathsf{iO}(C_0)) = 1] - \Pr_{(C_0, C_1) \leftarrow_{\mathcal{S}}} [\mathcal{D}(\mathsf{iO}(C_1)) = 1] \right| \leq \mathsf{negl}(n),$$

where the probabilities are over the randomness of the algorithms.

When obfuscating cryptographic algorithms with keys it is often convenient to use the notation $C[k](x)$ to denote the circuit with fixed encoded key(s) k and variable input x .

3 Hardware-Bound White-box Key Derivation Function

In this section, we first introduce our notion of a hardware module and explain how we instantiate it in our setting. Then, we provide the syntax and security notion for a hardware-bound key derivation function, present our construction and provide a security reduction to indistinguishability obfuscation and PPRFs.

3.1 Hardware Module

A schematic overview over the hardware module functionalities executed on the secure hardware is given in Figure 4. Namely, a hardware module comes with a key generation algorithm that generates the hardware master key k_{HWms} . This key generation algorithm is run at the manufacturer of the hardware and thus not depicted in Figure 4. The secure hardware allows to export a sub-key via querying the secure hardware with a *label*. The hardware then runs the *sub-key generation algorithm* $\mathsf{SubKgen}_{\mathsf{HW}}$ on k_{HWms} and *label* and returns the resulting sub-key k_{HWs1} . In addition, the secure hardware can be queried with a pair (x, label) . The hardware, then, uses the algorithm $\mathsf{Resp}_{\mathsf{HW}}$ to generate a PRF/MAC value σ for x under k_{HWms} . In order to avoid storing k_{HWs1} for all values *label*, the hardware re-derives k_{HWs1} anew each time $\mathsf{Resp}_{\mathsf{HW}}$ is called. The PRF/MAC value can be checked outside of the secure hardware by running $\mathsf{Check}_{\mathsf{SW}}$ on k_{HWs1} and the pair (x, σ) .

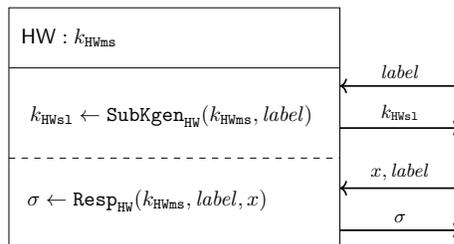


Fig. 4. Functionalities of the hardware module performed in the hardware. The **Check** operation is performed by the software program corresponding to the *label* value.

Remark. For secure use of the hardware-binding, it is important that the verification key k_{HWs1} is not stored in plain since, otherwise, one might derive PRF/MAC values using k_{HWs1} rather than querying the hardware. Thus, the Check_{SW} functionality will need to be white-boxed (and will later be bound to another functionality such as our white-box KDF), essentially making it asymmetric. Note that the syntax and correctness of our hardware module also allows to be directly implemented by a (standard, asymmetric) signature scheme such as provided by the Android Keystore [And18b]. In that case, the verification key does not need additional protection. We chose to implement the hardware with a symmetric primitive for efficiency and simplicity of the proof. Note that regardless of whether one uses MACs or signature schemes as a hardware functionality, the verification functionality needs to be cryptographically *bound to a software program* to be useful; else the software program can be code-lifted and run without performing the verification check.

Regardless of efficiency, both approaches, using signature schemes directly or making MAC verification asymmetric, are sound approaches. Importantly, in both cases, the soundness of the approach relies on *domain separation*, i.e., signatures/MAC for one *label* should not be mixed up with signatures/MACs for a different *label*. Finally, recall that fixed device identifiers (that do not depend on the input x) tend to provide very weak hardware-binding guarantees only, since once intercepted, they can be emulated for a code-lifted software program. We now give the syntax for our hardware-module.

Definition 11 (Hardware Module HWM). A hardware module HWM consists of four algorithms $(\text{Kgen}_{\text{HW}}, \text{SubKgen}_{\text{HW}}, \text{Resp}_{\text{HW}}, \text{Check}_{\text{SW}})$, where Kgen_{HW} and $\text{SubKgen}_{\text{HW}}$ are PPT algorithms with the following syntax

$$\begin{aligned} k_{\text{HWms}} &\leftarrow_{\text{s}} \text{Kgen}_{\text{HW}}(1^n) \\ k_{\text{HWs1}} &\leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWms}}, \text{label}) \end{aligned}$$

The algorithms Resp_{HW} and Check_{SW} are deterministic algorithms and have the following syntax

$$\begin{aligned} \sigma &\leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWms}}, \text{label}, x) \\ \{0, 1\} \ni b &\leftarrow \text{Check}_{\text{SW}}(k_{\text{HWs1}}, x, \sigma) \end{aligned}$$

Correctness requires that for all security parameters $n \in \mathbb{N}$,

$$\Pr[\text{Check}_{\text{SW}}(\text{SubKgen}_{\text{HW}}(k_{\text{HWms}}, \text{label}), x, \text{Resp}_{\text{HW}}(k_{\text{HWms}}, \text{label}, x)) = 1] = 1,$$

where the probability is over the sampling of $k_{\text{HWms}} \leftarrow_{\text{s}} \text{Kgen}_{\text{HW}}(1^n)$.

We do not define or prove security of the hardware module as a standalone primitive, since we will later prove security of our white-box KDF directly based on the puncturable PRFs used in the hardware module construction. Note that a standalone security definition for a hardware module could not capture MACs/PRFs and signature security simultaneously. In our hardware module

$\text{Exp}_{\text{WKDF}, \mathcal{A}}^{\text{IND-WKDF}}(1^n)$ <hr/> $b \leftarrow_{\$} \{0, 1\}$ $\mathcal{Q} := \emptyset$ $\text{label} \leftarrow_{\$} \mathcal{A}(1^n)$ $k_{\text{kdf}} \leftarrow_{\$} \text{Kgen}(1^n)$ $k_{\text{HWms}} \leftarrow_{\$} \text{Kgen}_{\text{HW}}(1^n)$ $k_{\text{HWs1}} \leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWms}}, \text{label})$ $\text{WKDF} \leftarrow_{\$} \text{Comp}(k_{\text{kdf}}, k_{\text{HWs1}})$ $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}_{\text{HW}}}^{\text{OResp, OKDF}}(\text{WKDF})$ $\text{return } b'$	$\text{OKDF}()$ <hr/> $e \leftarrow_{\$} \{0, 1\}^n$ $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ $\text{if } b = 0$ $\quad \hat{k} \leftarrow \text{KDF}(k_{\text{kdf}}, e)$ else $\quad \hat{k} \leftarrow_{\$} \{0, 1\}^n$ $\text{return } (e, \hat{k})$
$\text{OResp}(e)$ <hr/> $\text{assert } e \notin \mathcal{Q}$ $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ $\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWms}}, \text{label}, e)$ $\text{return } \sigma$	$\text{OSubKgen}(\text{label}')$ <hr/> $\text{assert } \text{label}' \neq \text{label}$ $k'_{\text{HWs1}} \leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWms}}, \text{label}')$ $\text{return } k'_{\text{HWs1}}$

Fig. 5. The $\text{Exp}_{\text{WKDF}, \mathcal{A}}^{\text{IND-WKDF}}(1^n)$ security game.

construction below, key generation samples a random key, sub-key generation applies a PRF to it to derive a sub-key, and Resp_{HW} and Check_{SW} essentially implement a PRF-based MAC. Note that the additional PRG evaluation in Check_{SW} is merely used to enable the proof technique by Sahai and Waters [SW14].

Construction 1. *Let PRG be a pseudorandom generator, PRF be a pseudorandom function and PPRF be a puncturable pseudorandom function. We construct a hardware module as follows:*

$\text{Kgen}_{\text{HW}}(1^n)$ <hr/> $1: k_{\text{HWms}} \leftarrow_{\$} \{0, 1\}^n$ $2: \text{return } k_{\text{HWms}}$	$\text{SubKgen}_{\text{HW}}(k_{\text{HWms}}, \text{label})$ <hr/> $1: k_{\text{HWs1}} \leftarrow \text{PRF}(k_{\text{HWms}}, \text{label})$ $2: \text{return } k_{\text{HWs1}}$
$\text{Resp}_{\text{HW}}(k_{\text{HWms}}, \text{label}, x)$ <hr/> $1: \sigma \leftarrow \text{PPRF}(\text{PRF}(k_{\text{HWms}}, \text{label}), x)$ $2: \text{return } \sigma$	$\text{Check}_{\text{SW}}(k_{\text{HWs1}}, x, \sigma)$ <hr/> $1: \text{if } \text{PRG}(\sigma) = \text{PRG}(\text{PPRF}(k_{\text{HWs1}}, x))$ $2: \text{return } 1 \text{ else return } 0$

Hardware-bound White-box Key Derivation Function. We now define and construct a hardware-bound white-box key derivation function WKDF. We here build on the previously introduced hardware module and a traditional KDF.

In the compiling phase, a compiler Comp takes as input the KDF key k_{kdf} and the sub-key k_{HWS1} for Check_{SW} . The compiler generates a program WKDF which takes as input a pair (e, σ) and, intuitively, first checks whether σ is valid for e under k_{HWS1} and, if so, evaluates the KDF on k_{kdf} and e . The role of the compiler, conceptually, is to return a program where the KDF operation is *bound* to the verification operation, i.e., the two functionalities cannot be separated from each other and the (outcome of the) verification cannot be manipulated.

Definition 12 (WKDF). *A white-box key derivation scheme with hardware binding WKDF consists of a hardware module HWM, a key derivation function KDF, and a PPT compiling algorithm Comp:*

$$\text{WKDF} \leftarrow_{\text{s}} \text{Comp}(k_{\text{kdf}}, k_{\text{HWS1}}).$$

For all genuine k_{HWS} , for all k_{kdf} , for all e , for $k_{\text{HWS1}} = \text{SubKgen}_{\text{HW}}(k_{\text{HWS}}, \text{label})$ and $\sigma = \text{Resp}_{\text{HW}}(k_{\text{HWS}}, \text{label}, e)$, we have $\Pr[\text{KDF}(k_{\text{kdf}}, e) = \text{WKDF}(e, \sigma)] = 1$ where the probability is taken over compiling $\text{WKDF} \leftarrow_{\text{s}} \text{Comp}(k_{\text{kdf}}, k_{\text{HWS1}})$.

Security Model. We now define security for a WKDF via the IND-WKDF security game, illustrated in Figure 5. We want to capture the pseudorandomness of keys derived from the WKDF, i.e., an adversary should not be able to distinguish between a key produced from the WKDF and an equally long key sampled at random. As a white-box adversary, the game provides the adversary with the capability of inspecting the WKDF itself (i.e. its circuit or implementation code). Recall that additionally, we want to capture the notion of hardware binding: if the adversary tries to run the WKDF without having access to its designated hardware, the WKDF should not be executable any more. We model this by giving the adversary hardware access via a OResp oracle, which the adversary queries by providing a context value e and which the adversary can query a polynomial number of times. With the reply σ from the hardware component, the adversary is able to run WKDF on the context value e used for querying OResp . Additionally, we grant the adversary access to an oracle OSubKgen . This oracle produces hardware sub-keys which might be used to generate new hardware-bound white-box key derivation functions. To avoid trivial attacks, the adversary is not allowed to request a sub-key under the same *label* that was used to generate the initial WKDF. To capture the pseudorandomness of derived keys, the adversary has access to a real-or-random oracle OKDF . The oracle OKDF first samples a new context value e and then, depending on a random bit b , it either returns the output of the KDF (under key k_{kdf}) or a random string of equal length. To avoid trivial attacks, the adversary is not allowed to query the OResp oracle with the same context value that was sampled by the challenger. Finally, \mathcal{A} outputs a bit b' and wins if $b' = b$.

Definition 13 (IND-WKDF). *We say that a hardware-bound white-box key derivation scheme WKDF is IND-WKDF-secure if all PPT adversaries \mathcal{A} have a negligible distinguishing advantage in game $\text{Exp}_{\text{WKDF}, \mathcal{A}}^{\text{IND-WKDF}}(1^n)$, see Figure 5.*

Note that one could also provide the adversary with a recompilation oracle [DLPR14] so that the adversary can request several (independent) copies of the WKDF based on the same key. We refrain from including this feature into our model, but note that our construction can be shown to achieve also this stronger notion, as indistinguishability obfuscation makes recompilation adversarially simulatable.

3.2 Construction of a WKDF

We now construct a WKDF, based on a traditional KDF and the previously introduced hardware module HWM. As discussed before, the compiler Comp , on input the KDF key k_{kdf} and the hardware sub-key k_{HWS1} binds the hardware Check_{SW} procedure to the KDF evaluation. Concretely, the compiler constructs a circuit $\text{C}[k_{\text{kdf}}, k_{\text{HWS1}}]$ and obfuscates it using indistinguishability obfuscation. The circuit $\text{C}[k_{\text{kdf}}, k_{\text{HWS1}}]$, on input (e, σ) first checks whether $\text{Check}_{\text{SW}}(k_{\text{HWS1}}, e, \sigma)$ equals 1. If yes, it returns the output of $\text{KDF}(k_{\text{kdf}}, e)$. Else, it returns the all-zero string. The reason that the construction is secure, is, intuitively, that the obfuscation of $\text{C}[k_{\text{kdf}}, k_{\text{HWS1}}]$ achieves the desired binding property. We now first give the KDF construction and then the WKDF construction directly below.

Construction 2. *Let PPRF be a puncturable pseudorandom function scheme, then we construct our KDF as follows:*

$\text{Kgen}_{\text{kdf}}(1^n)$	$\text{KDF}(k_{\text{kdf}}, e)$
1: $k_{\text{kdf}} \leftarrow_{\$} \{0, 1\}^n$	1: $\hat{k} \leftarrow \text{PPRF}(k_{\text{kdf}}, e)$
2: return k_{kdf}	2: return \hat{k}

Construction 3. *Let iO be an indistinguishability obfuscator. Based on the hardware module HWM given in Construction 1 and the key derivation scheme KDF given in Construction 2, we construct WKDF by defining the following compiler Comp :*

$\text{C}[k_{\text{kdf}}, k_{\text{HWS1}}](e, \sigma)$	$\text{Comp}(k_{\text{kdf}}, k_{\text{HWS1}})$
1: $v \leftarrow \text{Check}_{\text{SW}}(k_{\text{HWS1}}, e, \sigma)$	1: $\text{WKDF} \leftarrow_{\$} \text{iO}(\text{C}[k_{\text{kdf}}, k_{\text{HWS1}}])$
2: if $v = \perp$ return 0^n	2: return WKDF
3: else	
4: $\hat{k} \leftarrow \text{KDF}(k_{\text{kdf}}, e)$	
5: return \hat{k}	

Theorem 1. *Let PRG be a pseudorandom generator, let PRF be a pseudorandom function, let PPRF be a puncturable PRF, and iO be an indistinguishability obfuscator for appropriate p -admissible samplers. Then Construction 3 is a secure white-box KDF scheme WKDF.*

Proof. Let \mathcal{A} be a PPT adversary. Let $\text{Exp}_{\mathcal{A},b}^{\text{IND-WKDF}}$ denote the IND-WKDF game with a value $b \in \{0, 1\}$ hardcoded. We show that

$$\text{Exp}_{\mathcal{A},0}^{\text{IND-WKDF}}(1^n) \approx \text{Exp}_{\mathcal{A},1}^{\text{IND-WKDF}}(1^n).$$

Overview: The proof is a hybrid argument over the number of queries q that \mathcal{A} makes to the OKDF oracle which either returns the evaluation of a puncturable PRF or a random string of same length. Note that the proof applies the punctured programs approach by Sahai and Waters [SW14]. Our hybrid games maintain a counter j that increases by 1 whenever the adversary queries the OKDF oracle. In the i -th hybrid the challenger returns a random string whenever the counter $j > i$, otherwise it returns the evaluation of the PPRF. In other words, whenever we move to the next hybrid, the oracle returns an additional random string such that we sequentially replace PPRF values by random strings of appropriate size. After at most polynomial steps we have replaced all OKDF outputs by random values and obtain $\text{Exp}_{\mathcal{A},1}^{\text{IND-WKDF}}(1^n)$.

Detailed proof: Let $q(n)$ be a polynomial strict upper bound on the number of queries that \mathcal{A} makes to oracle OKDF. We define a sequence of adversary-dependent hybrid games Game_1^0 to $\text{Game}_1^{q(n)}$ such that

$$\text{Exp}_{\mathcal{A},0}^{\text{IND-WKDF}} \approx \text{Game}_1^0 \quad (1)$$

$$\text{Exp}_{\mathcal{A},1}^{\text{IND-WKDF}} \approx \text{Game}_1^{q(n)}. \quad (2)$$

18 game hops show that for $0 \leq i \leq q(n) - 1$:

$$\text{Game}_1^i \approx \text{Game}_{18}^i \quad (3) \quad \text{Game}_{18}^i \approx \text{Game}_1^{i+1}. \quad (4)$$

Indistinguishability of Game_1^0 and $\text{Game}_1^{q(n)}$ then follows by a standard hybrid argument, guessing the hybrid index at random. We define the games and specify the game-hops below. The circuits are depicted on the right of this page. Equation 4 follows by inspection of the definitions of Game_{18}^i and Game_1^{i+1} . We now turn to showing Equation 3 which is the technical heart of the theorem.

```

C1[kkdf, kHWS1](e, σ)
if PRG(σ) = PRG(PPRF(kHWS1, e))

    k̂ ← PPRF(kkdf, e)
    return k̂
else return 0n

C2[kkdf, z, kz, τ](e, σ)
if e = z and PRG(σ) = PRG(τ)
or if PRG(σ) = PRG(Eval(kz, e))
    k̂ ← PPRF(kkdf, e)
    return k̂
else return 0n

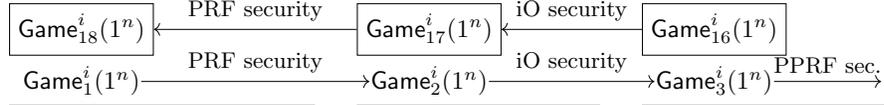
C3[kkdf, z, kz, y](e, σ)
if e = z and PRG(σ) = y
or if PRG(σ) = PRG(Eval(kz, e))

    k̂ ← PPRF(kkdf, e)
    return k̂
else return 0n

C4[kkdfz, z, kz, y, k](e, σ)
if e = z and PRG(σ) = y
    return k
if PRG(σ) = PRG(Eval(kz, e))
    k̂ ← Eval(kkdfz, e)
    return k̂
else return 0n

C5[kkdfz, z, kz, y](e, σ)
if e = z and PRG(σ) = y
    return 0n
if PRG(σ) = PRG(Eval(kz, e))
    k̂ ← Eval(kkdfz, e)
    return k̂
else return 0n

```



1: $z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ 2: $j \leftarrow 0$ 3: $label \leftarrow_{\$} \mathcal{A}(1^n)$ 4: $k_{kdf} \leftarrow_{\$} \{0, 1\}^n$ 5: 6: $k \leftarrow \text{PPRF}(k_{kdf}, z)$ 7: $k_{HWms} \leftarrow_{\$} \{0, 1\}^n$ 8: $k_{HWs1} \leftarrow_{\$} \text{PRF}(k_{HWms}, label)$ 9: 10: 11: $C \leftarrow C_1[k_{kdf}, k_{HWs1}]$ 12: $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ 13: $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResP}_{\text{HW}}, \text{OKDF}}(\text{WKDF})$ 14: return b'	1: $z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ 2: $j \leftarrow 0$ 3: $label \leftarrow_{\$} \mathcal{A}(1^n)$ 4: $k_{kdf} \leftarrow_{\$} \{0, 1\}^n$ 5: 6: $k \leftarrow \text{PPRF}(k_{kdf}, z)$ 7: $k_{HWms} \leftarrow_{\$} \{0, 1\}^n$ 8: $k_{HWs1} \leftarrow_{\$} \{0, 1\}^n$ 9: 10: 11: $C \leftarrow C_1[k_{kdf}, k_{HWs1}]$ 12: $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ 13: $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResP}_{\text{HW}}, \text{OKDF}}(\text{WKDF})$ 14: return b'	1: $z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ 2: $j \leftarrow 0$ 3: $label \leftarrow_{\$} \mathcal{A}(1^n)$ 4: $k_{kdf} \leftarrow_{\$} \{0, 1\}^n$ 5: 6: $k \leftarrow \text{PPRF}(k_{kdf}, z)$ 7: $k_{HWms} \leftarrow_{\$} \{0, 1\}^n$ 8: $k_{HWs1} \leftarrow_{\$} \{0, 1\}^n$ 9: $k_z \leftarrow_{\$} \text{Punct}(k_{HWs1}, z)$ 10: $\tau \leftarrow \text{PPRF}(k_{HWs1}, z)$ 11: $C \leftarrow C_2[k_{kdf}, z, k_z, \tau]$ 12: $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ 13: $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResP}_{\text{HW}}, \text{OKDF}}(\text{WKDF})$ 14: return b'
---	--	---

OResp(e)

1: **assert** $e \notin \mathcal{Q}$
2: $\mathcal{Q} := \mathcal{Q} \cup \{e\}$
3: $y \leftarrow \text{PPRF}(k_{HWs1}, e, 1^n)$
4: **return** y

OResp(e)

assert $e \notin \mathcal{Q}$
 $\mathcal{Q} := \mathcal{Q} \cup \{e\}$
 $y \leftarrow \text{PPRF}(k_{HWs1}, e, 1^n)$
return y

OResp(e)

assert $e \notin \mathcal{Q}$
 $\mathcal{Q} := \mathcal{Q} \cup \{e\}$
 $\tilde{y} \leftarrow \text{Eval}(k_z, e, 1^n)$
return y

OKDF()

1: $j \leftarrow j + 1$
2: **if** $i = j$ $k' \leftarrow_{\$} \{0, 1\}^n$
3: **return** (z, k) (z, k')
4: **else** $e \leftarrow_{\$} \{0, 1\}^n$
5: $\mathcal{Q} := \mathcal{Q} \cup \{e\}$
6: **if** $j > i$
7: $\hat{k} \leftarrow \text{PPRF}(k_{kdf}, e)$
8: **else** $\hat{k} \leftarrow_{\$} \{0, 1\}^n$
9: **return** (e, \hat{k})

OKDF()

$j \leftarrow j + 1$
if $i = j$ $k' \leftarrow_{\$} \{0, 1\}^n$
return (z, k) (z, k')
else $e \leftarrow_{\$} \{0, 1\}^n$
 $\mathcal{Q} := \mathcal{Q} \cup \{e\}$
if $j > i$
 $\hat{k} \leftarrow \text{PPRF}(k_{kdf}, e)$
else $\hat{k} \leftarrow_{\$} \{0, 1\}^n$
return (e, \hat{k})

OKDF()

$j \leftarrow j + 1$
if $i = j$ $k' \leftarrow_{\$} \{0, 1\}^n$
return (z, k) (z, k')
else $e \leftarrow_{\$} \{0, 1\}^n$
 $\mathcal{Q} := \mathcal{Q} \cup \{e\}$
if $j > i$
 $\hat{k} \leftarrow \text{PPRF}(k_{kdf}, e)$
else $\hat{k} \leftarrow_{\$} \{0, 1\}^n$
return \hat{k}

OSubKgen_{HW}($label'$)

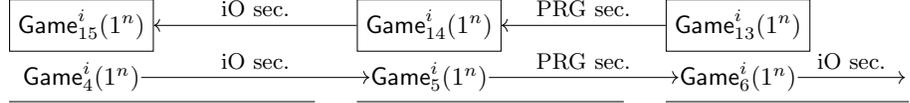
1: **assert** $label' \neq label$
2: $k'_{HWs1} \leftarrow \text{PRF}(k_{HWms}, label')$
3: **return** k'_{HWs1}

OSubKgen_{HW}($label'$)

assert $label' \neq label$
 $k'_{HWs1} \leftarrow \text{PRF}(k_{HWms}, label')$
return k'_{HWs1}

OSubKgen_{HW}($label'$)

assert $label' \neq label$
 $k'_{HWs1} \leftarrow \text{PRF}(k_{HWms}, label')$
return k'_{HWs1}



```

1:  $z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ 
2:  $j \leftarrow 0$ 
3:  $label \leftarrow_{\$} \mathcal{A}(1^n)$ 
4:  $k_{\text{kdf}} \leftarrow_{\$} \{0, 1\}^n$ 
5:
6:  $k \leftarrow \text{PPRF}(k_{\text{kdf}}, z)$ 
7:  $k_{\text{HWms}} \leftarrow_{\$} \{0, 1\}^n$ 
8:  $k_{\text{HWs1}} \leftarrow_{\$} \{0, 1\}^n$ 
9:  $k_z \leftarrow_{\$} \text{Punct}(k_{\text{HWs1}}, z)$ 
10:  $\tau \leftarrow_{\$} \{0, 1\}^n$ 
11:  $C \leftarrow C_2[k_{\text{kdf}}, z, k_z, \tau]$ 
12:  $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ 
13:  $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResPiv}, \text{OKDF}}(\text{WKDF})$ 
14: return  $b'$ 

```

OResp(e)

```

1: assert  $e \notin \mathcal{Q}$ 
2:  $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
3:  $y \leftarrow \text{Eval}(k_z, e, 1^n)$ 
4: return  $y$ 

```

OKDF()

```

1:  $j \leftarrow j + 1$ 
2: if  $i = j$   $k' \leftarrow_{\$} \{0, 1\}^n$ 
3: return  $(z, k)$   $(z, k')$ 
4: else  $e \leftarrow_{\$} \{0, 1\}^n$ 
5:  $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
6: if  $j > i$ 
7:  $\hat{k} \leftarrow \text{PPRF}(k_{\text{kdf}}, e)$ 
8: else  $\hat{k} \leftarrow_{\$} \{0, 1\}^n$ 
9: return  $\hat{k}$ 

```

OSubKgen_{HW}($label'$)

```

1: assert  $label' \neq label$ 
2:  $k'_{\text{HWs1}} \leftarrow \text{PRF}(k_{\text{HWms}}, label')$ 
3: return  $k'_{\text{HWs1}}$ 

```

```

1:  $z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ 
2:  $j \leftarrow 0$ 
3:  $label \leftarrow_{\$} \mathcal{A}(1^n)$ 
4:  $k_{\text{kdf}} \leftarrow_{\$} \{0, 1\}^n$ 
5:
6:  $k \leftarrow \text{PPRF}(k_{\text{kdf}}, z)$ 
7:  $k_{\text{HWms}} \leftarrow_{\$} \{0, 1\}^n$ 
8:  $k_{\text{HWs1}} \leftarrow_{\$} \{0, 1\}^n$ 
9:  $k_z \leftarrow_{\$} \text{Punct}(k_{\text{HWs1}}, z)$ 
10:  $\tau \leftarrow_{\$} \{0, 1\}^n; y \leftarrow \text{PRG}(\tau)$ 
11:  $C \leftarrow C_3[k_{\text{kdf}}, z, k_z, y]$ 
12:  $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ 
13:  $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResPiv}, \text{OKDF}}(\text{WKDF})$ 
14: return  $b'$ 

```

OResp(e)

```

1: assert  $e \notin \mathcal{Q}$ 
2:  $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
3:  $y \leftarrow \text{Eval}(k_z, e, 1^n)$ 
4: return  $y$ 

```

OKDF()

```

1:  $j \leftarrow j + 1$ 
2: if  $i = j$   $k' \leftarrow_{\$} \{0, 1\}^n$ 
3: return  $(z, k)$   $(z, k')$ 
4: else  $e \leftarrow_{\$} \{0, 1\}^n$ 
5:  $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
6: if  $j > i$ 
7:  $\hat{k} \leftarrow \text{PPRF}(k_{\text{kdf}}, e)$ 
8: else  $\hat{k} \leftarrow_{\$} \{0, 1\}^n$ 
9: return  $\hat{k}$ 

```

OSubKgen_{HW}($label'$)

```

1: assert  $label' \neq label$ 
2:  $k'_{\text{HWs1}} \leftarrow \text{PRF}(k_{\text{HWms}}, label')$ 
3: return  $k'_{\text{HWs1}}$ 

```

```

1:  $z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ 
2:  $j \leftarrow 0$ 
3:  $label \leftarrow_{\$} \mathcal{A}(1^n)$ 
4:  $k_{\text{kdf}} \leftarrow_{\$} \{0, 1\}^n$ 
5:
6:  $k \leftarrow \text{PPRF}(k_{\text{kdf}}, z)$ 
7:  $k_{\text{HWms}} \leftarrow_{\$} \{0, 1\}^n$ 
8:  $k_{\text{HWs1}} \leftarrow_{\$} \{0, 1\}^n$ 
9:  $k_z \leftarrow_{\$} \text{Punct}(k_{\text{HWs1}}, z)$ 
10:  $y \leftarrow_{\$} \{0, 1\}^{2n}$ 
11:  $C \leftarrow C_3[k_{\text{kdf}}, z, k_z, y]$ 
12:  $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ 
13:  $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResPiv}, \text{OKDF}}(\text{WKDF})$ 
14: return  $b'$ 

```

OResp(e)

```

1: assert  $e \notin \mathcal{Q}$ 
2:  $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
3:  $y \leftarrow \text{Eval}(k_z, e, 1^n)$ 
4: return  $y$ 

```

OKDF()

```

1:  $j \leftarrow j + 1$ 
2: if  $i = j$   $k' \leftarrow_{\$} \{0, 1\}^n$ 
3: return  $(z, k)$   $(z, k')$ 
4: else  $e \leftarrow_{\$} \{0, 1\}^n$ 
5:  $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
6: if  $j > i$ 
7:  $\hat{k} \leftarrow \text{PPRF}(k_{\text{kdf}}, e)$ 
8: else  $\hat{k} \leftarrow_{\$} \{0, 1\}^n$ 
9: return  $\hat{k}$ 

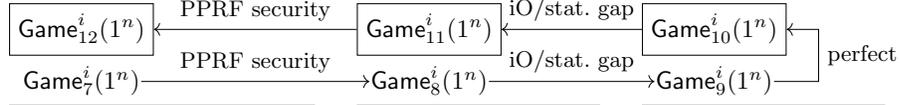
```

OSubKgen_{HW}($label'$)

```

1: assert  $label' \neq label$ 
2:  $k'_{\text{HWs1}} \leftarrow \text{PRF}(k_{\text{HWms}}, label')$ 
3: return  $k'_{\text{HWs1}}$ 

```



1: $z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ 2: $j \leftarrow 0$ 3: $label \leftarrow_{\$} \mathcal{A}(1^n)$ 4: $k_{kdf} \leftarrow_{\$} \{0, 1\}^n$ 5: $k_{kdf}^z \leftarrow_{\$} \text{Punct}(k_{kdf}, z)$ 6: $k \leftarrow \text{PPRF}(k_{kdf}, z)$ 7: $k_{HWms} \leftarrow_{\$} \{0, 1\}^n$ 8: $k_{HWs1} \leftarrow_{\$} \{0, 1\}^n$ 9: $k_z \leftarrow_{\$} \text{Punct}(k_{HWs1}, z)$ 10: $y \leftarrow_{\$} \{0, 1\}^{2n}$ 11: $C \leftarrow C_4[k_{kdf}^z, z, k_z, y, k]$ 12: $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ 13: $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResP}_{\text{HW}}, \text{OKDF}}(\text{WKDF})$ 14: return b'	$z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ $j \leftarrow 0$ $label \leftarrow_{\$} \mathcal{A}(1^n)$ $k_{kdf} \leftarrow_{\$} \{0, 1\}^n$ $k_{kdf}^z \leftarrow_{\$} \text{Punct}(k_{kdf}, z)$ $k \leftarrow_{\$} \{0, 1\}^n$ $k_{HWms} \leftarrow_{\$} \{0, 1\}^n$ $k_{HWs1} \leftarrow_{\$} \{0, 1\}^n$ $k_z \leftarrow_{\$} \text{Punct}(k_{HWs1}, z)$ $y \leftarrow_{\$} \{0, 1\}^{2n}$ $C \leftarrow C_4[k_{kdf}^z, z, k_z, y, k]$ $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResP}_{\text{HW}}, \text{OKDF}}(\text{WKDF})$ return b'	$z \leftarrow_{\$} \{0, 1\}^n, \mathcal{Q} := \{z\}$ $j \leftarrow 0$ $label \leftarrow_{\$} \mathcal{A}(1^n)$ $k_{kdf} \leftarrow_{\$} \{0, 1\}^n$ $k_{kdf}^z \leftarrow_{\$} \text{Punct}(k_{kdf}, z)$ $k \leftarrow_{\$} \{0, 1\}^n$ $k_{HWms} \leftarrow_{\$} \{0, 1\}^n$ $k_{HWs1} \leftarrow_{\$} \{0, 1\}^n$ $k_z \leftarrow_{\$} \text{Punct}(k_{HWs1}, z)$ $y \leftarrow_{\$} \{0, 1\}^{2n}$ $C \leftarrow C_5[k_{kdf}^z, z, k_z, y]$ $\text{WKDF} \leftarrow_{\$} \text{iO}(C, 1^n)$ $b' \leftarrow_{\$} \mathcal{A}_{\text{OSubKgen}}^{\text{ResP}_{\text{HW}}, \text{OKDF}}(\text{WKDF})$ return b'
---	--	---

OResp(e)

```

1: assert  $e \notin \mathcal{Q}$ 
2:  $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
3:  $y \leftarrow \text{Eval}(k_z, e, 1^n)$ 
4: return  $y$ 

```

OResp(e)

```

assert  $e \notin \mathcal{Q}$ 
 $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
 $y \leftarrow \text{Eval}(k_z, e, 1^n)$ 
return  $y$ 

```

OResp(e)

```

assert  $e \notin \mathcal{Q}$ 
 $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
 $y \leftarrow \text{Eval}(k_z, e, 1^n)$ 
return  $y$ 

```

OKDF()

```

1:  $j \leftarrow j + 1$ 
2: if  $i = j$   $k' \leftarrow_{\$} \{0, 1\}^n$ 
3:   return  $(z, k)$   $(z, k')$ 
4: else  $e \leftarrow_{\$} \{0, 1\}^n$ 
5:    $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
6:   if  $j > i$ 
7:      $\hat{k} \leftarrow \text{Eval}(k_{kdf}^z, e)$ 
8:   else  $\hat{k} \leftarrow_{\$} \{0, 1\}^n$ 
9:   return  $(e, \hat{k})$ 

```

OKDF()

```

 $j \leftarrow j + 1$ 
if  $i = j$   $k' \leftarrow_{\$} \{0, 1\}^n$ 
   return  $(z, k)$   $(z, k')$ 
else  $e \leftarrow_{\$} \{0, 1\}^n$ 
    $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
   if  $j > i$ 
      $\hat{k} \leftarrow \text{Eval}(k_{kdf}^z, e)$ 
   else  $\hat{k} \leftarrow_{\$} \{0, 1\}^n$ 
   return  $(e, \hat{k})$ 

```

OKDF()

```

 $j \leftarrow j + 1$ 
if  $i = j$   $k' \leftarrow_{\$} \{0, 1\}^n$ 
   return  $(z, k)$   $(z, k')$ 
else  $e \leftarrow_{\$} \{0, 1\}^n$ 
    $\mathcal{Q} := \mathcal{Q} \cup \{e\}$ 
   if  $j > i$ 
      $\hat{k} \leftarrow \text{Eval}(k_{kdf}^z, e)$ 
   else  $\hat{k} \leftarrow_{\$} \{0, 1\}^n$ 
   return  $(e, \hat{k})$ 

```

OSubKgen_{HW}($label'$)

```

1: assert  $label' \neq label$ 
2:  $k'_{HWs1} \leftarrow \text{PRF}(k_{HWms}, label')$ 
3: return  $k'_{HWs1}$ 

```

OSubKgen_{HW}($label'$)

```

assert  $label' \neq label$ 
 $k'_{HWs1} \leftarrow \text{PRF}(k_{HWms}, label')$ 
return  $k'_{HWs1}$ 

```

OSubKgen_{HW}($label'$)

```

assert  $label' \neq label$ 
 $k'_{HWs1} \leftarrow \text{PRF}(k_{HWms}, label')$ 
return  $k'_{HWs1}$ 

```

We now reduce each game-hop to the underlying assumption. We omit boilerplate code for simulations and focus instead on describing the conceptual reasons why the reductions work. We only discuss the game-hops in the forward direction, since reductions in the backward direction proceed analogously.

Game 1 to Game 2. This game hop reduces to the PRF security of the PRF keyed with k_{Hwms} . We here rely on the fact that $\text{OSubKgen}_{\text{HW}}(\text{label}')$ does not allow to query $\text{PRF}(k_{\text{Hwms}}, \cdot)$ on label .

Game 2 to Game 3. This game hop reduces to iO security and relies on the correctness of the PPRF. The correctness of the PPRF implies that the first two lines of \mathbf{C}_1 and \mathbf{C}_2 are equivalent, which allows to apply iO security. Note that the oracle OResp cannot be queried on z since z is added to \mathcal{Q} in the very beginning of the game. Therefore, it suffices to use the punctured key k_z in OResp .

Game 3 to Game 4. This game hop reduces to the IND-PPRF security of PPRF keyed with k_{Hws1} and punctured at z . For the reduction, it is important to note that throughout the game, only punctured versions k_z of k_{Hws1} are used, except for calculating τ .

Game 4 to Game 5. This game hop replaces $\mathbf{C}_2[k_{\text{kdf}}, z, k_z, \tau]$ by $\mathbf{C}_3[k_{\text{kdf}}, z, k_z, y]$ and relies on iO security. Instead of hardcoding τ into \mathbf{C}_2 and computing y as $\text{PRG}(\tau)$ within circuit \mathbf{C}_2 , the value $y = \text{PRG}(\tau)$ is directly hardcoded into circuit \mathbf{C}_3 . As the two circuits are functionally equivalent, the game hop reduces to iO security.

Game 5 to Game 6. This game hop replaces $y = \text{PRG}(\tau)$ by a randomly sampled value y . This game hop reduces to PRG security, since the variable τ that is sampled in $\text{Game}_5^i(1^n)$ is not used anywhere else in the game(s).

Game 6 to Game 7. This game hop reduces to iO security with an additional negligible statistical loss. In detail, the uniformity of the sampling of z from $\{0, 1\}^n$ ensures that, with overwhelming probability, the oracle OKDF does not return $(e, *)$ for a different than the i th query of \mathcal{A} to the OKDF oracle, i.e., the change to OKDF yields a negligible statistical difference between $\text{Game}_6^i(1^n)$ and $\text{Game}_7^i(1^n)$. The more important change is the use of k_{kdf}^z in \mathbf{C}_4 . Due to the correctness of the puncturable PRF, the circuits \mathbf{C}_3 and \mathbf{C}_4 are functionally equivalent and thus, this game hop can be reduced to iO security.

Game 7 to Game 8. This game hop reduces to the IND-PPRF security of PPRF, keyed with k_{kdf} and punctured at z . For the reduction, it is important to note that throughout the game, only the punctured version k_{kdf}^z of k_{kdf} is used, except for calculating k .

Game 8 to Game 9. Note that with overwhelming probability, y is not in the image of the PRG, since the image of the PRG is of size 2^n only, whereas y is sampled from a set of size 2^{2n} . Therefore, y is most likely outside the image of the PRG. If it is, then the circuits C_4 and C_5 are functionally equivalent as the **if** condition in the first line of C_4 cannot be satisfied by any input. Thus, this game hop reduces to iO security.

Game 9 to Game 10. Importantly, C_5 does not depend on k anymore and thus, it is perfectly indistinguishable for the adversary whether the OKDF uses k or an independently drawn value k' that OKDF samples in the moment of the i th query.

Game 10 to Game 18. These game hops are analogous in the backward direction. Note that k is not used in OKDF anymore in the game hops from Game 10 to Game 18.

Connecting the hybrids. We now show that Equation 4 holds, which we recall is

$$\text{Game}_{18}^i \approx \text{Game}_1^{i+1}$$

On a high-level, Game_{18}^i and Game_1^{i+1} are identical since they both sample the first $i + 1$ keys in the OKDF oracle at random and compute the remaining keys using the PPRF. Note that Game_{18}^i and Game_1^{i+1} are only identical up to a negligible statistical difference since the pre-sampled value z is consumed at query i to the OKDF oracle in Game_{18}^i and only in query $i + 1$ in Game_1^{i+1} . However, as z is sampled uniformly at random from $\{0, 1\}^n$, z remains statistically hidden from the adversary until it is returned as an output from the OKDF oracle. Thus, it is infeasible to determine when the pre-sampled value z was returned. We omit a detailed code-comparison, since it is quite simple.

Connecting the hybrids to the original game. Finally, we show Equation 1 and Equation 2. We start with the former which we recall is

$$\text{Exp}_{\mathcal{A},0}^{\text{IND-WKDF}} \approx \text{Game}_1^0.$$

On a high-level, $\text{Exp}_{\mathcal{A},0}^{\text{IND-WKDF}}$ and Game_1^0 are identical since they both compute all keys in the OKDF oracle using the PPRF. Moreover, due to the correctness of the hardware module, it is functionally equivalent to use $\text{PPRF}(k_{\text{HWS1}}, e, 1^n)$ and $\text{PPRF}(\text{PRF}(k_{\text{HWS1}}, e, 1^n))$ in oracle OResp. However, there is a negligible statistical difference between $\text{Exp}_{\mathcal{A},0}^{\text{IND-WKDF}}$ and Game_1^0 , since Game_1^0 pre-samples a uniformly random value z from $\{0, 1\}^n$, while $\text{Exp}_{\mathcal{A},0}^{\text{IND-WKDF}}$ does not. However, as z is sampled uniformly at random, it remains statistically hidden from the adversary until it is returned as an output from the OKDF oracle. We omit a detailed code-comparison, since it is quite simple. The reasoning for Equation 2 is analogous.

4 Secure Payment Application

In this section we build a secure payment scheme, assuming an IND-WKDF-secure WKDF, as constructed in the previous section. As mentioned before, the main idea of the construction of the payment application is to bind a symmetric encryption scheme (that is only known to provide black-box security) on top of the WKDF via a layer of indistinguishability obfuscation and thereby bootstrap the security of the WKDF to the security of the symmetric encryption scheme and the entire payment application. We first introduce the syntax of a payment application, then present our construction and then turn to the security reduction for the construction.

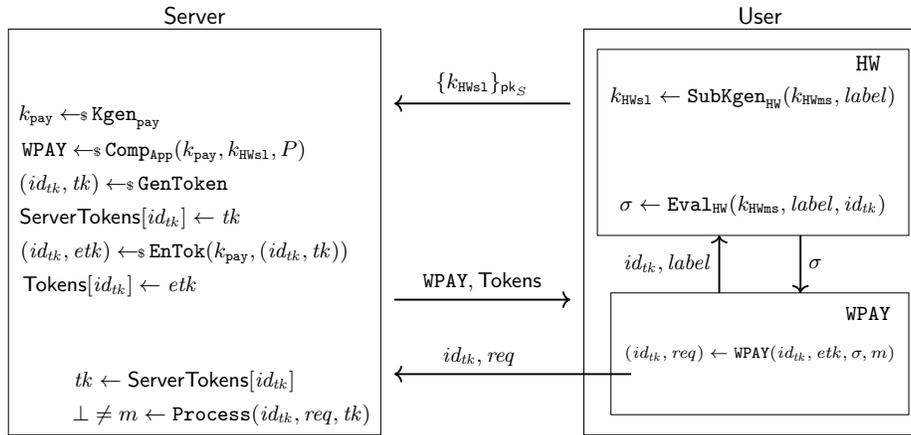


Fig. 6. Diagram of our payment scheme. The hardware calculates the value σ via the Eval_{HW} function on input id_{tk} , the label and the master key. WPAY executes the Check_{SW} function on input σ and on the sub-key and id_{tk} .

We now describe the process flow of an abstract payment application, illustrated in Figure 6. Note that this abstract payment application might also be implementable differently than assuming a WKDF. A payment application relies on a hardware module (see Definition 11), which we recall has a master hardware key k_{HWms} , allows to derive sub-keys $k_{\text{HWs1}} \leftarrow \text{SubKgen}(k_{\text{HWms}}, label)$ from the master hardware key and allows to obtain MAC/PRF values $\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWms}}, label, x)$ that can be verified outside of the hardware via the algorithm $\text{Check}_{\text{SW}}(k_{\text{HWs1}}, x, \sigma)$ that returns 0 or 1.

As the payment application is bound to a hardware module, the user starts by deriving a hardware sub-key $k_{\text{HWs1}} \leftarrow \text{SubKgen}(k_{\text{HWms}}, label)$ in their hardware and transmit it securely to the server. In Figure 6, we hint at this secure transmission of k_{HWs1} via an encryption under the server's public-key pk_S . In our model, we

later refrain from modeling this off-band transmission of k_{HWS1} and simply assume that it is implemented securely.

The server then draws a symmetric key k_{pay} for the user and binds it to the user’s hardware sub-key k_{HWS1} via the compilation algorithm Comp_{App} :

$$\text{WPAY} \leftarrow_{\text{s}} \text{Comp}_{\text{App}}(k_{\text{pay}}, k_{\text{HWS1}}, P).$$

The predicate function P restricts the set of valid messages that can be encrypted via WPAY . An example for useful restrictions are limits on the amount of the payment or hardcoding of the user’s payment data. Note that potentially, P can also contain cryptographic functionalities (which we do not model).

As we consider a tokenized payment scheme, in addition to WPAY the server also generates several tokens, encrypts them under k_{pay} and stores the encrypted tokens in an array Tokens that is indexed by token identifiers. It then transmits WPAY to the user, together with the array Tokens of encrypted tokens, see Figure 6.

Now, the user can use WPAY to generate requests to the server. To do so, WPAY takes as input a message m as well as a pair of a token identifier id_{tk} and its corresponding encrypted token etk . Conceptually, the goal of WPAY is to return a request req to the server that contains an encryption of m under the unencrypted token contained in etk . To facilitate verification on the server’s side, the user’s WPAY will also return the token identifier id_{tk} :

$$(id_{tk}, req) \leftarrow \text{WPAY}(id_{tk}, etk, \sigma, m)$$

Importantly, in addition to the aforementioned inputs, WPAY also takes as input σ , which is a hardware value obtained from making an Resp_{HW} query to the hardware for $(label, id_{tk})$. To ensure the hardware-binding, conceptually, WPAY needs to evaluate the algorithm $\text{Check}_{\text{SW}}(k_{\text{HWS1}}, id_{tk}, \sigma)$ internally and only perform the desired operations based on this check succeeding. Intuitively, the Check_{SW} operation also needs to be bound to all further operations of WPAY .

Finally, upon receiving (id_{tk}, req) , the server retrieves the token tk corresponding to id_{tk} and processes the request req via the algorithm Process . If the request is accepted, Process returns the message m that the client encrypted. Else, Process returns an error symbol \perp .

Note that we require the server to know the secret key k_{pay} to encrypt the tokens under k_{pay} before sending it to the user. The advantage of this design is that the values of the tokens are not exposed before being stored. Note that we consider the server to be a trusted and secure party which is a necessary assumption: As the server is in charge of generating the tokens, the server knows the token values anyway. Thus, the server additionally knowing the value of the secret key k_{kdf} of the user does not compromise the security of the mobile payment application from the perspective of the user.

Definition 14 (Hardware-Bound White-Box Payment Scheme).

A hardware-bound white-box payment scheme WPAY is parametrized with a message length parameter $\ell(n)$ and consists of a hardware module HWM and the following algorithms:

- $k_{\text{pay}} \leftarrow_{\text{s}} \text{Kgen}_{\text{pay}}(1^n)$: This randomized algorithm takes as input the security parameter and outputs the secret payment key k_{pay} ;
- $(id_{tk}, tk) \leftarrow_{\text{s}} \text{GenToken}(1^n)$: This randomized algorithm takes as input the security parameter and returns a token tk and a token identifier id_{tk} , where we assume that identifiers are unique with overwhelming probability;
- $(id_{tk}, etk) \leftarrow_{\text{s}} \text{EnTok}(k_{\text{pay}}, (id_{tk}, tk))$: This randomized algorithm takes as input a token tk together with its corresponding identifier id_{tk} and outputs an encrypted token etk , together with its corresponding token identifier id_{tk} ;
- $\text{WPAY} \leftarrow_{\text{s}} \text{CompApp}(k_{\text{pay}}, k_{\text{HWS1}}, P)$: This randomized algorithm takes as input a payment key k_{pay} , a hardware-binding key k_{HWS1} and a message filtering predicate $P : \{0, 1\}^{\ell(n)} \rightarrow \{0, 1\}$. It outputs a white-box payment application WPAY with syntax $(id_{tk}, req) \leftarrow \text{WPAY}(id_{tk}, etk, \sigma, m)$, where σ denotes a hardware value σ , $m \in \{0, 1\}^{\ell(n)}$ denotes a payment message, and req constitutes a payment request;
- $m \leftarrow \text{Process}(id_{tk}, req, tk)$: This deterministic algorithm takes as input a token identifier id_{tk} , a token value tk and a request req . It outputs a message m or \perp .

Moreover, we require that the following correctness property holds: For all keys k_{pay} , for all master hardware keys k_{HWS1} , for all pairs of tokens and token identifier (id_{tk}, tk) , for all predicates P , for all messages $m \in \{0, 1\}^{\ell(n)}$ such that $P(m) = 1$, for $\sigma = \text{Resp}_{\text{HW}}(k_{\text{HWS1}}, \text{label}, id_{tk})$, it holds that

$$\Pr[\text{Process}(\text{WPAY}(id_{tk}, etk, \sigma, m), tk) = m] = 1,$$

where the probability is taken over compiling $\text{WPAY} \leftarrow_{\text{s}} \text{CompApp}(k_{\text{pay}}, k_{\text{HWS1}}, P)$ and encrypting the token $(id_{tk}, etk) \leftarrow_{\text{s}} \text{EnTok}(k_{\text{pay}}, (id_{tk}, tk))$.

4.1 Security of White-Box Payment Applications

We now specify security of a white-box payment application scheme WPAY . Correctness of WPAY ensures that when having access to the hardware, the application WPAY is useful to generate payment requests. In turn, hardware-binding security ensures that when not having access to the hardware, then the application becomes useless. In other words, in absence of the hardware, the adversary cannot generate new requests and does not learn anything about the content of the requests sent to the server. Thus, the desired security properties in the absence of the hardware are the following:

- (1) Integrity of the requests transmitted from user to server.
- (2) Confidentiality of the messages contained in the requests transmitted from user to server.

We capture both properties via the IND-WPAY security game, depicted in Figure 7. IND-WPAY starts with a setup phase where the relevant keys are sampled, first for the hardware (line 3 and 4) and then for the payment application (line 5). Then, WPAY is compiled (line 6). Note that we allow the adversary to

$\text{Exp}_{\text{WPAY}, \mathcal{A}}^{\text{IND-WPAY}}(1^n)$	$\text{OTransaction}(m_0, m_1)$
<ol style="list-style-type: none"> 1 : $b \leftarrow_{\\$} \{0, 1\}$ 2 : $(label, P) \leftarrow_{\\$} \mathcal{A}(1^n)$ 3 : $k_{\text{HWms}} \leftarrow_{\\$} \text{Kgen}_{\text{HW}}(1^n)$ 4 : $k_{\text{HWs1}} \leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWms}}, label)$ 5 : $k_{\text{pay}} \leftarrow_{\\$} \text{Kgen}_{\text{pay}}(1^n)$ 6 : $\text{WPAY} \leftarrow_{\\$} \text{Comp}_{\text{App}}(k_{\text{pay}}, k_{\text{HWs1}}, P)$ 7 : $b^* \leftarrow_{\\$} \mathcal{A}_{\text{OProcess, OTransaction}}^{\text{OResp, OSubKgen, OGetTok}}(\text{WPAY})$ 8 : return $(b^* = b)$ 	<ol style="list-style-type: none"> 1 : assert $m_0 = m_1 = \ell(n)$ 2 : assert $P(m_0) = 1 \wedge P(m_1) = 1$ 3 : $(id_{tk}, tk) \leftarrow_{\\$} \text{GenToken}(1^n)$ 4 : $\text{ServerTokens}[id_{tk}] \leftarrow tk$ 5 : $(id_{tk}, etk) \leftarrow_{\\$} \text{EnTok}(k_{\text{pay}}, (id_{tk}, tk))$ 6 : $\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWms}}, label, id_{tk})$ 7 : $(id_{tk}, req) \leftarrow \text{WPAY}(id_{tk}, etk, \sigma, m_b)$ 8 : $C := C \cup \{req\}$ 9 : return (id_{tk}, etk, req)
<hr style="border: 0.5px solid black;"/> $\text{OResp}()$ <ol style="list-style-type: none"> 1 : $(id_{tk}, tk) \leftarrow_{\\$} \text{GenToken}(1^n)$ 2 : $\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWms}}, label, id_{tk})$ 3 : $\text{RespUsed}[id_{tk}] \leftarrow 1$ 4 : $(id_{tk}, etk) \leftarrow_{\\$} \text{EnTok}(k_{\text{pay}}, (id_{tk}, tk))$ 5 : return (id_{tk}, etk, σ) 	<hr style="border: 0.5px solid black;"/> $\text{OProcess}(id_{tk}, req)$ <ol style="list-style-type: none"> 1 : assert $\text{RespUsed}[id_{tk}] = 0$ 2 : assert $req \notin C$ 3 : $tk \leftarrow \text{ServerTokens}[id_{tk}]$ 4 : $m \leftarrow \text{Process}(id_{tk}, req, tk)$ 5 : $\text{ServerTokens}[id_{tk}] \leftarrow \perp$ 6 : if $b = 0$ return m 7 : else $b = 1$ return \perp
<hr style="border: 0.5px solid black;"/> $\text{OSubKgen}(label')$ <ol style="list-style-type: none"> 1 : assert $label' \neq label$ 2 : $k'_{\text{HWs1}} \leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWms}}, label')$ 3 : return k'_{HWs1} 	
<hr style="border: 0.5px solid black;"/> $\text{OGetTok}()$ <ol style="list-style-type: none"> 1 : $(id_{tk}, tk) \leftarrow_{\\$} \text{GenToken}(1^n)$ 2 : $\text{ServerTokens}[id_{tk}] \leftarrow tk$ 3 : $(id_{tk}, etk) \leftarrow_{\\$} \text{EnTok}(k_{\text{pay}}, (id_{tk}, tk))$ 4 : return (id_{tk}, etk) 	

Fig. 7. $\text{Exp}_{\text{WPAY}, \mathcal{A}}^{\text{IND-WPAY}}(1^n)$ game capturing integrity and confidentiality.

choose the filter function P , modelling that security should hold for all possible filter functions. We also allow the adversary to choose the hardware *label*. In practice, neither P nor *label* are adversarially chosen, but giving this ability to the adversary in the setup phase only makes our model stronger. Note that we consider our adversary as stateful, i.e., the adversary in the setup phase (line 2) shares state with the adversary that accesses oracles (line 7) in order to find out the secret bit b . As the adversary is a white-box adversary, it receives the compiled **WPAY** as input (line 7). We now turn to the explanation of the oracles.

Oracles **OResp** and **OSubKgen** model the adversary’s hardware access. Upon querying **OResp**, a pair of token and its respective identifier are sampled at random. The value of the token identifier is used for generating the hardware value σ , and the oracle returns all three values. The adversary is thus able to run their own **WPAY** to generate a request message for that specific token, token id and hardware value.

OSubKgen lets the adversary observe hardware sub-key values. Such values could be values used by different applications or values from older versions of a payment application. To avoid trivial attacks, the adversary is not allowed to obtain the sub-key k_{HWS1} that was used for compiling the white-box application **WPAY** (line 6 in the main $\text{Exp}_{\text{WPAY}, \mathcal{A}}^{\text{IND-WPAY}}$ procedure).

Oracle **OGetTok** models generation and encryption of the tokens on the server’s side and storing them in a list. We recall that a white-box adversary might be able to access the encrypted tokens stored on a user’s phone and thus, **OGetTok** returns the encrypted token and its identifier to the adversary.

We now turn to the two remaining oracles that encode the desired security properties of confidentiality and integrity. We start with the transaction oracle **OTransaction**, which encodes the confidentiality property. On a conceptual level, **OTransaction** plays a similar role as the left-or-right encryption oracle in the security game for authenticated encryption (c.f. Definition 2): It encrypts either the left or the right message, depending on whether the secret bit b is 0 or 1. Upon submitting the two messages to the **OTransaction** oracle, the oracle randomly samples a pair of a token and its respective identifier and the token value is encrypted via the **EnTok** algorithm. The oracle then generates the necessary hardware value σ based on the sampled token identifier, and then generates a request message based on the token identifier, the encrypted token, the value σ and one of the two messages submitted by the adversary. The oracle then saves the generated request on a list C and returns the token identifier, the encrypted token and the request message.

Note that the adversary is not able to use their own app to generate the same request message received by the transaction oracle. Namely, while the adversary is able to choose which message to encrypt with their own app, the token used for encrypting the message is chosen at random and only with negligible probability will both request messages look the same.

The process oracle **OProcess** encodes the integrity of the request messages similarly to the decryption oracle in authenticated encryption. The adversary can submit arbitrary values as long as those were not obtained from **OTransaction**

(check if they are in the set C) or if they were not generated using an id_{tk} generated by the `OResponse` oracle for generating the hardware value σ . The server retrieves (line 3) the token tk corresponding to id_{tk} from the token list, decrypts (line 4) the request req with id_{tk} and tk using the `Process` algorithm and deletes tk from the token list. Authenticity is modeled by only returning the message to the adversary if $b = 0$ and returning an error if $b = 1$. Thereby, the adversary is able to learn the secret bit b whenever the adversary is able to forge a fresh request.

Remark. Note that many useful properties are implied by our security definition. For instance, IND-WPAY security implies that the token values remain secret, unless the adversary queries the hardware on id_{tk} .

Definition 15 (IND-WPAY Security). *A hardware-bound white-box payment application scheme WPAY is said to be IND-WPAY-secure if all PPT adversaries \mathcal{A} have negligible distinguishing advantage in the game $\text{Exp}_{\text{WPAY}, \mathcal{A}}^{\text{IND-WPAY}}(1^n)$, specified in Figure 7.*

4.2 Construction of White-box Payment Scheme

We now construct a white-box payment scheme WPAY, which is IND-WPAY-secure (see Figure 7), assuming the IND-WKDF-security of a white-box key derivation function WKDF (see Figure 5). We first give an overview of the algorithms of our construction. First, `Kgenpay` randomly samples a payment key k_{pay} and the `GenToken` algorithm randomly samples a token tk and its respective identifier id_{tk} . `EnTok` encrypts a token tk in the following way: For each token with the identifier id_{tk} , it generates a key $\hat{k} \leftarrow \text{KDF}(k_{\text{pay}}, id_{tk})$, and uses an authenticated encryption scheme to encrypt the token using \hat{k} . I.e., each token is encrypted using a different key. Note that each key is generated based on the same key k_{pay} , but based on a different context value id_{tk} since each token has a unique identifier.

The compilation algorithm `CompApp` of the payment scheme takes as input a payment key k_{pay} , a sub-key value k_{HWS1} as well as a message filtering predicate P . It uses the compilation algorithm `Comp` of the WKDF and runs it on $(k_{\text{pay}}, k_{\text{HWS1}})$ to obtain a hardware-bound white-box program WKDF. It then runs indistinguishability obfuscation on a circuit $\text{C}[\text{WKDF}, P]$ and returns the output of the obfuscation as WPAY. The circuit $\text{C}[\text{WKDF}, P]$ and WPAY provide the same functionality, but in WPAY, due to the layer of obfuscation, the different operations cannot be separated anymore. $\text{C}[\text{WKDF}, P]$ takes as input a pair (id_{tk}, etk) of an encrypted token and its token identifier as well as a hardware value σ and a message m . It first runs WKDF on (id_{tk}, σ) to obtain an output \hat{k} , and recall that the security of WKDF ensures that \hat{k} is only a KDF (and not an error value) if σ is the correct hardware value that yields $\text{Check}_{\text{SW}}(k_{\text{HWS1}}, id_{tk}, \sigma) = 1$. Thus, the hardware-binding of WPAY is directly inherited from the hardware-binding of WKDF. Now, $\text{C}[\text{WKDF}, P]$ uses \hat{k} to decrypt the etk (line 3), checks whether $P(m) = 1$ (line 4) and, if so, encrypts the message m using the token tk as key

$\text{Kgen}_{\text{pay}}(1^n)$	$\text{C}[\text{WKDF}, P](id_{tk}, etk, \sigma, m)$	$\text{GenToken}(1^n)$
1 : $k_{\text{kdf}} \leftarrow_{\$} \text{Kgen}(1^n)$	1 : $e \leftarrow id_{tk}$	1 : $tk \leftarrow_{\$} \{0, 1\}^n$
2 : $k_{\text{pay}} \leftarrow k_{\text{kdf}}$	2 : $nc \leftarrow id_{tk}$	2 : $id_{tk} \leftarrow_{\$} \{0, 1\}^n$
3 : return k_{pay}	3 : $c_1 \leftarrow etk$	3 : return (id_{tk}, tk)
$\text{Comp}_{\text{App}}(k_{\text{pay}}, k_{\text{HWS1}}, P)$	4 : $\hat{k} \leftarrow \text{WKDF}(e, \sigma)$	$\text{EnTok}(k_{\text{pay}}, (id_{tk}, tk))$
1 : $\text{WKDF} \leftarrow_{\$} \text{Comp}(k_{\text{pay}}, k_{\text{HWS1}})$	5 : $tk \leftarrow \text{Dec}_1(\hat{k}, c_1, nc)$	1 : $nc \leftarrow id_{tk}$
2 : $\text{WPAY} \leftarrow_{\$} \text{iO}(\text{C}[\text{WKDF}, P])$	6 : if $ m = \ell(n)$ and $P(m) = 1$	2 : $e \leftarrow id_{tk}$
3 : return WPAY	7 : $c_2 \leftarrow \text{Enc}_2(tk, m, nc)$	3 : $\hat{k} \leftarrow \text{KDF}(k_{\text{pay}}, e)$
$\text{Process}(id_{tk}, req, tk)$	8 : $req \leftarrow c_2$	4 : $c_1 \leftarrow \text{Enc}_1(\hat{k}, tk, nc)$
1 : $nc \leftarrow id_{tk}$	9 : else $req \leftarrow \perp$	5 : $etk \leftarrow c_1$
2 : $c_2 \leftarrow req$	10 : return (id_{tk}, req)	6 : return (id_{tk}, etk)
3 : $m' \leftarrow \text{Dec}_2(tk, c_2, nc)$		
4 : return m'		

Fig. 8. Construction 4 of a WPAY scheme

(line 5-6) and returns the resulting ciphertext as a request message req (line 8). Note that in the construction, we need to avoid the use of randomness, since we cannot rely on the randomness being honestly generated. Thus, we use a nonce-based encryption scheme, and the nonce nc used for encryption of the token is not only retrieved (line 2) for decryption of the token (line 3), but also re-used for encrypting the request (line 6). Note that the nonce in etk is not malleable since we encrypt the tokens with an authenticated encryption scheme which provides ciphertext integrity.

Recall that WPAY returns not only req but also the token identifier id_{tk} of the token that was used to encrypt m . This is because the use of tk authenticates the user towards the server which, on its side, retrieves the token tk corresponding to id_{tk} , and runs the Process algorithm on (req, tk) which decrypts req with tk and returns the result.

Construction 4. *Based on two authenticated encryption schemes $(\text{Enc}_1, \text{Dec}_1)$, $(\text{Enc}_2, \text{Dec}_2)$ and the WKDF in Construction 3, we construct a white-box payment scheme with hardware binding $\text{WPAY} = (\text{Kgen}_{\text{pay}}, \text{GenToken}, \text{EnTok}, \text{Comp}_{\text{App}}, \text{Process})$, see Figure 4.2.*

We provide a proof of the following theorem in Appendix A.

Theorem 2. *Let $(\text{Enc}_1, \text{Dec}_1)$ and $(\text{Enc}_2, \text{Dec}_2)$ be two AE-secure symmetric encryption schemes, let WKDF be a IND-WKDF-secure white-box key derivation scheme, and let iO be an indistinguishability obfuscator for appropriate p -admissible samplers. Then the white-box payment scheme WPAY in Construction 4 is IND-WPAY-secure.*

Acknowledgments

Marc Fischlin and Christian Janson have been [co-]funded by the DFG as part of projects S4 and P2 within the CRC 1119 CROSSING.

References

- ABB⁺19. Gildas Avoine, Muhammed Ali Bingöl, Ioana Boureanu, Srdjan Capkun, Gerhard P. Hancke, Süleyman Kardas, Chong Hee Kim, Cédric Lauradoux, Benjamin Martin, Jorge Munilla, Alberto Peinado, Kasper Bonne Rasmussen, Dave Singelée, Aslan Tchamkerten, Rolando Trujillo-Rasua, and Serge Vaudenay. Security of distance-bounding: A survey. *ACM Comput. Surv.*, 51(5):94:1–94:33, 2019.
- ABBB⁺19. Estuardo Alpirez Bock, Joppe W. Bos, Chris Brzuska, Charles Hubain, Wil Michiels, Cristofaro Mune, Eloi Sanfelix Gonzalez, Philippe Teuwen, and Alexander Treff. White-box cryptography: Don’t forget about grey-box attacks. *Journal of Cryptology*, Feb 2019.
- AHKM14. Daniel Apon, Yan Huang, Jonathan Katz, and Alex J. Malozemoff. Implementing cryptographic program obfuscation. Cryptology ePrint Archive, Report 2014/779, 2014. <http://eprint.iacr.org/2014/779>.
- And18a. Android Developers. Build. Class Documentation, Last retrieved: October 2018. <https://developer.android.com/reference/android/os/Build>.
- And18b. Android Developers. Keystore. Class Documentation, Last retrieved: October 2018. <https://developer.android.com/reference/java/security/KeyStore>.
- And18c. Android Developers. Telephony manager. Class Documentation, Last retrieved: October 2018. <https://developer.android.com/reference/android/telephony/TelephonyManager>.
- BAB⁺19. Estuardo Alpirez Bock, Alessandro Amadori, Joppe W. Bos, Chris Brzuska, and Wil Michiels. Doubly half-injective PRGs for incompressible white-box cryptography. In *CT-RSA 2019*, LNCS, pages 189–209. Springer, Heidelberg, 2019.
- BBIJ17. Subhadeep Banik, Andrey Bogdanov, Takanori Isobe, and Martin Bjerregaard Jepsen. Analysis of software countermeasures for whitebox encryption. Cryptology ePrint Archive, Report 2017/183, 2017. <http://eprint.iacr.org/2017/183>.
- BBK14. Alex Biryukov, Charles Bouillaguet, and Dmitry Khovratovich. Cryptographic schemes based on the ASASA structure: Black-box, white-box, and public-key (extended abstract). In Palash Sarkar and Tetsu Iwata, editors, *ASIACRYPT 2014, Part I*, volume 8873 of *LNCS*, pages 63–84. Springer, Heidelberg, December 2014.
- BI15. Andrey Bogdanov and Takanori Isobe. White-box cryptography revisited: Space-hard ciphers. In Indrajit Ray, Ninghui Li, and Christopher Kruegel, editors, *ACM CCS 15*, pages 1058–1069. ACM Press, October 2015.
- BIT16. Andrey Bogdanov, Takanori Isobe, and Elmar Tischhauser. Towards practical whitebox cryptography: Optimizing efficiency and space hardness. In Jung Hee Cheon and Tsuyoshi Takagi, editors, *ASIACRYPT 2016, Part I*, volume 10031 of *LNCS*, pages 126–158. Springer, Heidelberg, December 2016.

- BKTS16. Thomas Bocek, Christian Killer, Christos Tsias, and Burkhard Stiller. An nfc relay attack with off-the-shelf hardware and software. In Rémi Badonnel, Robert Koch, Aiko Pras, Martin Drašar, and Burkhard Stiller, editors, *Management and Security in the Age of Hyperconnectivity*, pages 71–83, Cham, 2016. Springer International Publishing.
- BN00. Mihir Bellare and Chanathip Namprempre. Authenticated encryption: Relations among notions and analysis of the generic composition paradigm. In Tatsuaki Okamoto, editor, *ASIACRYPT 2000*, volume 1976 of *LNCS*, pages 531–545. Springer, Heidelberg, December 2000.
- BSW16. Mihir Bellare, Igors Stepanovs, and Brent Waters. New negative results on differing-inputs obfuscation. In Marc Fischlin and Jean-Sébastien Coron, editors, *EUROCRYPT 2016, Part II*, volume 9666 of *LNCS*, pages 792–821. Springer, Heidelberg, May 2016.
- BW13. Dan Boneh and Brent Waters. Constrained pseudorandom functions and their applications. In Kazuo Sako and Palash Sarkar, editors, *ASIACRYPT 2013, Part II*, volume 8270 of *LNCS*, pages 280–300. Springer, Heidelberg, December 2013.
- CdRP14. Tim Cooijmans, Joeri de Ruiter, and Erik Poll. Analysis of secure key storage solutions on android. In *Proceedings of the 4th ACM Workshop on Security and Privacy in Smartphones & Mobile Devices*, SPSM '14, pages 11–20. ACM, 2014.
- CEJv03. Stanley Chow, Philip A. Eisen, Harold Johnson, and Paul C. van Oorschot. White-box cryptography and an AES implementation. In Kaisa Nyberg and Howard M. Heys, editors, *SAC 2002*, volume 2595 of *LNCS*, pages 250–270. Springer, Heidelberg, August 2003.
- CEJv003. Stanley Chow, Philip A. Eisen, Harold Johnson, and Paul C. van Oorschot. A white-box DES implementation for DRM applications. In Joan Feigenbaum, editor, *Security and Privacy in Digital Rights Management, ACM CCS-9 Workshop, DRM 2002*, volume 2696 of *LNCS*, pages 1–15. Springer, 2003.
- cyb19. cybercrypt. Ches 2019 capture the flag challenge - the whibox contest - edition 2, 2019. <https://www.cyber-crypt.com/whibox-contest/>.
- DLPR14. Cécile Delerablée, Tancrede Lepoint, Pascal Paillier, and Matthieu Rivain. White-box security notions for symmetric encryption schemes. In Tanja Lange, Kristin Lauter, and Petr Lisonek, editors, *SAC 2013*, volume 8282 of *LNCS*, pages 247–264. Springer, Heidelberg, August 2014.
- ECR17. ECRYPT. Ches 2017 capture the flag challenge - the whibox contest, 2017. <https://whibox.cr.yt.to/>.
- FHMM11. Lishoy Francis, Gerhard Hancke, Keith Mayes, and Konstantinos Markantonakis. Practical relay attack on contactless transactions by using NFC mobile phones. *Cryptology ePrint Archive*, Report 2011/618, 2011. <http://eprint.iacr.org/2011/618>.
- FKKM16. Pierre-Alain Fouque, Pierre Karpman, Paul Kirchner, and Brice Minaud. Efficient and provable white-box primitives. In Jung Hee Cheon and Tsuyoshi Takagi, editors, *ASIACRYPT 2016, Part I*, volume 10031 of *LNCS*, pages 159–188. Springer, Heidelberg, December 2016.
- GPRW18. Louis Goubin, Pascal Paillier, Matthieu Rivain, and Junwei Wang. How to reveal the secrets of an obscure white-box implementation. *Cryptology ePrint Archive*, Report 2018/098, 2018. <https://eprint.iacr.org/2018/098>.

- Han05. Gerhard Hancke. A practical relay attack on iso 14443 proximity cards. Technical report, 2005.
- Kra10. Hugo Krawczyk. Cryptographic extraction and key derivation: The HKDF scheme. In Tal Rabin, editor, *CRYPTO 2010*, volume 6223 of *LNCS*, pages 631–648. Springer, Heidelberg, August 2010.
- Mas17. Mastercard. Mastercard mobile payment sdk, 2017. <https://developer.mastercard.com/media/32/b3/b6a8b4134e50bfe53590c128085e/mastercard-mobile-payment-sdk-security-guide-v2.0.pdf>.
- Riv17. Matthieu Rivain. White-box cryptography. Presentation CARDIS 2017, 2017. <http://www.matthieurivain.com/files/slides-cardis17.pdf>.
- Rog02. Phillip Rogaway. Authenticated-encryption with associated-data. In Vijayalakshmi Atluri, editor, *ACM CCS 02*, pages 98–107. ACM Press, November 2002.
- Rog04. Phillip Rogaway. Nonce-based symmetric encryption. In Bimal K. Roy and Willi Meier, editors, *FSE 2004*, volume 3017 of *LNCS*, pages 348–359. Springer, Heidelberg, February 2004.
- SdHM15. Eloi Sanfelix, Job de Haas, and Cristofaro Mune. Unboxing the white-box: Practical attacks against obfuscated ciphers. Presentation at BlackHat Europe 2015, 2015. <https://www.blackhat.com/eu-15/briefings.html>.
- Sma14. Smart Card Alliance Mobile and NFC Council. Host card emulation 101. white paper, 2014. <http://www.smartcardalliance.org/downloads/HCE-101-WP-FINAL-081114-clean.pdf>.
- SW14. Amit Sahai and Brent Waters. How to use indistinguishability obfuscation: deniable encryption, and more. In David B. Shmoys, editor, *46th ACM STOC*, pages 475–484. ACM Press, May / June 2014.
- Wys12. Brecht Wyseur. White-box cryptography: Hiding keys in software. 2012. <http://www.whiteboxcrypto.com/research.php>.

A Proof of Theorem 2

Let \mathcal{A} be a PPT adversary. Let $\text{Exp}_{\mathcal{A},0}^{\text{IND-WPAY}}(1^n)$ denote the IND-WPAY game with $b = 0$ hardcoded and let $\text{Exp}_{\mathcal{A},1}^{\text{IND-WPAY}}(1^n)$ denote the IND-WPAY game with $b = 1$ hardcoded. We need to show that \mathcal{A} has negligible distinguishing advantage, i.e., that the probability that \mathcal{A} returns 1 in $\text{Exp}_{\mathcal{A},0}^{\text{IND-WPAY}}(1^n)$ differs from the probability that \mathcal{A} returns 1 in $\text{Exp}_{\mathcal{A},1}^{\text{IND-WPAY}}(1^n)$ at most by a negligible amount.

For the following proof overview, observe that in the IND-WPAY game, only the `OProcess` oracle and the `OTransaction` oracle depend on the bit b . Namely, `OTransaction` encrypts m_b , and `OProcess` implements ideal authentication if $b = 1$ by simply rejecting all adversarially generated requests.

On a high-level, the security proof proceeds as follows. (1) We replace all keys generated by the `WKDF` with random keys in `OTransaction` and `OGetTok` and reduce this game hop to the security of the `WKDF`. (2) In the `OTransaction` and `OGetTok` oracles, instead of encrypting the tokens tk with SE_1 , we encrypt $0^{|tk|}$. To do so, we make a hybrid argument over the number of queries that the adversary makes to `OTransaction` and `OGetTok` and for each such query make a reduction to the AE security of SE_1 . (Note that IND-CPA security would

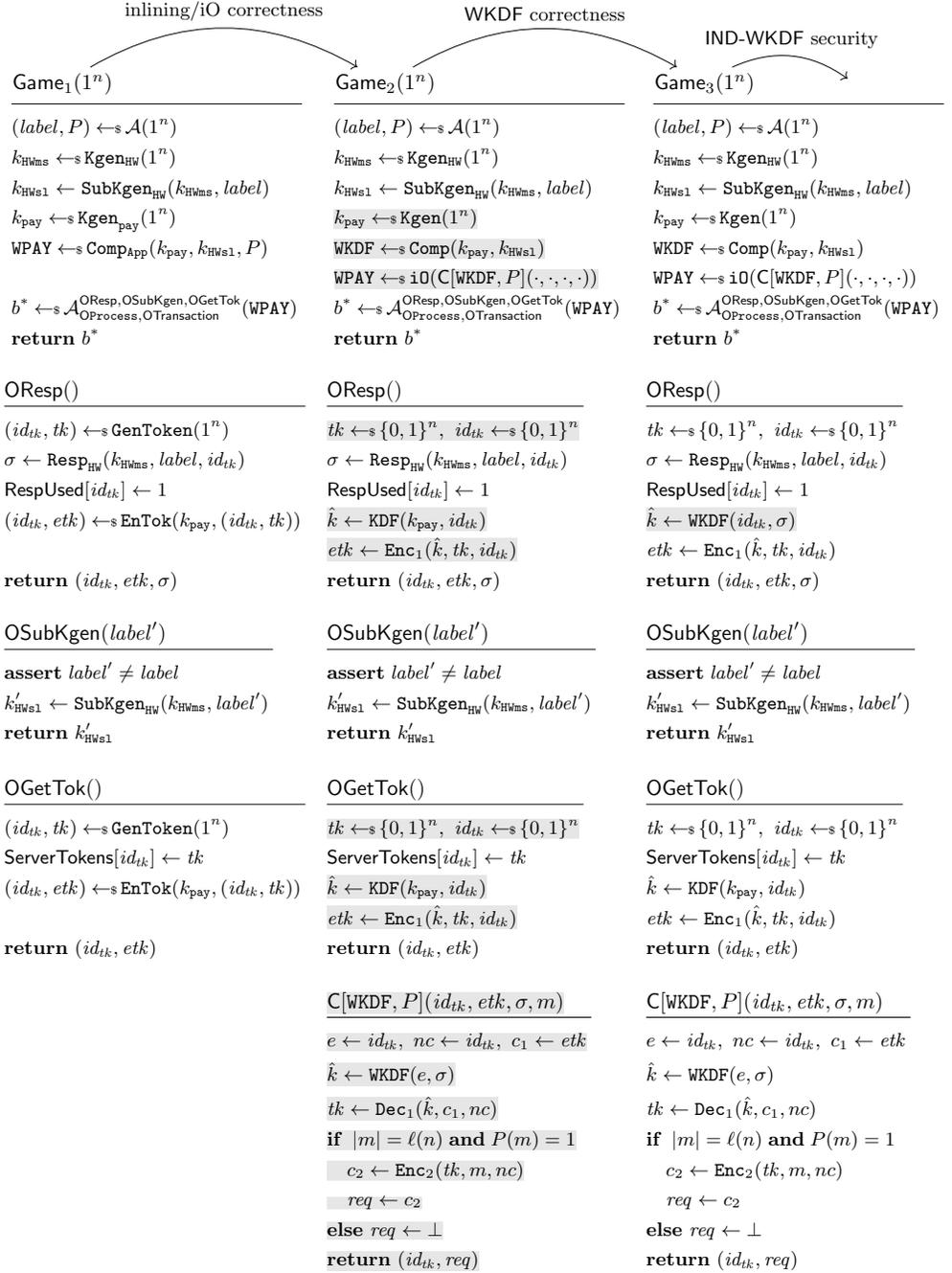
suffice here, since the game provides no decryption oracle for SE_1). (3) In the `OTransaction` oracle, instead of encrypting m_0 , we now encrypt m_1 . At the same time, in the `OProcess` oracle, we do not perform encryptions anymore but rather answer all adversarial queries by \perp . This proof again proceeds via a hybrid argument over the number of queries to `OTransaction` and `OGetTok`, since one token value is generated in each of these calls and we need to reduce to the security to each of them. (4) We now de-idealize SE_1 again and encrypt the real token values. (5) We de-idealize the `WKDF`.

Proof. We now provide formal descriptions of the game hops implementing step (1), (2) and (3). Step (4) is analogous to step (2), and step (5) is analogous to step (1). Below, Game_1 is $\text{Exp}_{\text{WPAY}, \mathcal{A}}^{\text{IND-WPAY}}(1^n)$ with bit b hardcoded to $b = 0$. Note that due to the length of the game, the description is split over two pages. Also recall that we measure the probability that the adversary returns 0, i.e., instead of returning $(b = b^*)$, the game only returns b^* . This is possible by a standard transformation on advantage terms with a loss of a factor of 2. Game_2 is $\text{Exp}_{\text{WPAY}, \mathcal{A}}^{\text{IND-WPAY}}(1^n)$ with Construction 4 inlined and `OTransaction` uses the code of $C[\text{WKDF}, P]$ directly rather than using the corresponding obfuscated version. This replacement is justified, as the indistinguishability obfuscator preserves perfect functional equivalence between the obfuscated and the unobfuscated program. From Game_2 to Game_3 , in the `OResp`, we use `WKDF` instead of `KDF` which is justified by the correctness of `WKDF`. Moreover, in the `OTransaction` oracle, we do not recalculate \hat{k} , as the correctness of `WKDF` ensures that `WKDF` and `KDF` compute the same key.

Game₃ to Game₄ We need to reduce this game hop to the IND-WKDF-security of the `WKDF`. The reduction receives (P) from the adversary \mathcal{A} and passes on *label* to the IND-KDF game. It then receives `WKDF` from the IND-WKDF game and computes $\text{WPAY} \leftarrow \text{iO}(C[\text{WKDF}, P])$ and returns `WPAY` to \mathcal{A} and then simulates its oracles and returns the same bit as \mathcal{A} . To simulate the `OGetTok` and the `OTransaction` oracle, the reduction uses its `OKDF()` oracle. As all oracles sample the input to the key derivation function uniformly at random, the simulation is sound. Moreover, the reduction uses the `OResp` oracle of the IND-WKDF game to simulate the `OResp` oracle of the IND-WPAY game. Note that the **assert** condition is, with statistically overwhelming probability not violated, because the inputs are chosen at random by the reduction to simulate the `OResp` oracle of the IND-WPAY game. The reduction knows the token values and, thus, can simulate `OProcess` perfectly. Finally, the reduction uses `OSubKgen` to simulate `OSubKgen` by simply forwarding queries.

Game₄ to Game₅ We need to reduce this game hop to the AE-security of SE_1 . We carry out a hybrid argument over the number of queries to `OTransaction` and `OGetTok`. Note that this is possible, because these oracles use uniformly random values to key SE_1 and these values are not used at other points in the system.

Game₅ to Game₆ We need to reduce this game hop to the AE-security of SE_2 . We carry out a hybrid argument over the number of queries to $OTransaction$ and $OGetTok$. Note that this is possible, because these oracles use uniformly random values to key SE_2 and these values are only used in $OTransaction$, $OGetTok$ and $OProcess$. Second, note that the nonces used for each encryption are randomly generated token identifiers id_{tk} . Therefore, the restriction given in line 4 of the $Transaction$ oracle in the $IND-WPAY$ -game corresponds to the second line of the encryption oracle in the AE game. This means that an adversary in the $IND-WPAY$ has the same restrictions for encrypting a message as an adversary in the AE game. Additionally, the process oracle only accepts request messages that were generated via the $Transaction$ oracle. This restriction aligns with the restriction of the decryption oracle in the AE game.



OTransaction(m_0, m_1)

assert $|m_0| = |m_1| = \ell(n)$
assert $P(m_0) = 1 \wedge P(m_1) = 1$
 $(id_{tk}, tk) \leftarrow_s \text{GenToken}(1^n)$
 $\text{ServerTokens}[id_{tk}] \leftarrow tk$
 $(id_{tk}, etk) \leftarrow_s \text{EnTok}(k_{\text{pay}}, (id_{tk}, tk))$

 $\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{Hwms}}, label, id_{tk})$
 $(id_{tk}, req) \leftarrow \text{WPAY}(id_{tk}, etk, \sigma, m_0)$

$C := C \cup \{req\}$
return (id_{tk}, etk, req)

OProcess(id_{tk}, req)

assert $\text{RespUsed}[id_{tk}] = 0$
assert $req \notin C$
 $tk \leftarrow \text{ServerTokens}[id_{tk}]$
 $m \leftarrow \text{Process}(id_{tk}, req, tk)$
 $\text{ServerTokens}[id_{tk}] \leftarrow \perp$
return m

OTransaction(m_0, m_1)

assert $|m_0| = |m_1| = \ell(n)$
assert $P(m_0) = 1 \wedge P(m_1) = 1$
 $tk \leftarrow_s \{0, 1\}^n, id_{tk} \leftarrow_s \{0, 1\}^n$
 $\text{ServerTokens}[id_{tk}] \leftarrow tk$
 $\hat{k} \leftarrow \text{KDF}(k_{\text{pay}}, id_{tk})$
 $etk \leftarrow \text{Enc}_1(\hat{k}, tk, id_{tk})$
 $\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{Hwms}}, label, id_{tk})$
 $\hat{k} \leftarrow \text{WKDF}(id_{tk}, \sigma)$
 $tk \leftarrow \text{Dec}_1(\hat{k}, etk, id_{tk})$
if $|m| = \ell(n)$ **and** $P(m) = 1$
 $req \leftarrow \text{Enc}_2(tk, m_0, nc)$
else $req \leftarrow \perp$
 $C := C \cup \{req\}$
return (id_{tk}, etk, req)

OProcess(id_{tk}, req)

assert $\text{RespUsed}[id_{tk}] = 0$
assert $req \notin C$
 $tk \leftarrow \text{ServerTokens}[id_{tk}]$
 $\hat{m} \leftarrow \text{Dec}_2(tk, req, id_{tk})$
 $\text{ServerTokens}[id_{tk}] \leftarrow \perp$
return m

OTransaction(m_0, m_1)

assert $|m_0| = |m_1| = \ell(n)$
assert $P(m_0) = 1 \wedge P(m_1) = 1$
 $tk \leftarrow_s \{0, 1\}^n, id_{tk} \leftarrow_s \{0, 1\}^n$
 $\text{ServerTokens}[id_{tk}] \leftarrow tk$
 $\hat{k} \leftarrow \text{KDF}(k_{\text{pay}}, id_{tk})$
 $etk \leftarrow \text{Enc}_1(\hat{k}, tk, id_{tk})$

$tk \leftarrow \text{Dec}_1(\hat{k}, etk, id_{tk})$
if $|m| = \ell(n)$ **and** $P(m) = 1$
 $req \leftarrow \text{Enc}_2(tk, m_0, nc)$
else $req \leftarrow \perp$
 $C := C \cup \{req\}$
return (id_{tk}, etk, req)

OProcess(id_{tk}, req)

assert $\text{RespUsed}[id_{tk}] = 0$
assert $req \notin C$
 $tk \leftarrow \text{ServerTokens}[id_{tk}]$
 $m \leftarrow \text{Dec}_2(tk, req, id_{tk})$
 $\text{ServerTokens}[id_{tk}] \leftarrow \perp$
return m



OTransaction(m_0, m_1)

assert $|m_0| = |m_1| = \ell(n)$
assert $P(m_0) = 1 \wedge P(m_1) = 1$
 $tk \leftarrow_s \{0, 1\}^n$, $id_{tk} \leftarrow_s \{0, 1\}^n$
 $\text{ServerTokens}[id_{tk}] \leftarrow tk$
 $\hat{k} \leftarrow \{0, 1\}^\lambda$
 $etk \leftarrow \text{Enc}_1(\hat{k}, tk, id_{tk})$

$tk \leftarrow \text{Dec}_1(\hat{k}, etk, id_{tk})$
if $|m| = \ell(n)$ **and** $P(m) = 1$
 $req \leftarrow \text{Enc}_2(tk, m_0, nc)$
else $req \leftarrow \perp$
 $C := C \cup \{req\}$
return (id_{tk}, etk, req)

OProcess(id_{tk}, req)

assert $\text{RespUsed}[id_{tk}] = 0$
assert $req \notin C$
 $tk \leftarrow \text{ServerTokens}[id_{tk}]$
 $m \leftarrow \text{Dec}_2(tk, req, id_{tk})$
 $\text{ServerTokens}[id_{tk}] \leftarrow \perp$
return m

OTransaction(m_0, m_1)

assert $|m_0| = |m_1| = \ell(n)$
assert $P(m_0) = 1 \wedge P(m_1) = 1$
 $tk \leftarrow_s \{0, 1\}^n$, $id_{tk} \leftarrow_s \{0, 1\}^n$
 $\text{ServerTokens}[id_{tk}] \leftarrow tk$
 $\hat{k} \leftarrow \{0, 1\}^\lambda$
 $etk \leftarrow \text{Enc}_1(\hat{k}, 0^{|tk|}, id_{tk})$

$tk \leftarrow \text{Dec}_1(\hat{k}, etk, id_{tk})$
if $|m| = \ell(n)$ **and** $P(m) = 1$
 $req \leftarrow \text{Enc}_2(tk, m_0, nc)$
else $req \leftarrow \perp$
 $C := C \cup \{req\}$
return (id_{tk}, etk, req)

OProcess(id_{tk}, req)

assert $\text{RespUsed}[id_{tk}] = 0$
assert $req \notin C$
 $tk \leftarrow \text{ServerTokens}[id_{tk}]$
 $m \leftarrow \text{Dec}_2(tk, req, id_{tk})$
 $\text{ServerTokens}[id_{tk}] \leftarrow \perp$
return m

OTransaction(m_0, m_1)

assert $|m_0| = |m_1| = \ell(n)$
assert $P(m_0) = 1 \wedge P(m_1) = 1$
 $tk \leftarrow_s \{0, 1\}^n$, $id_{tk} \leftarrow_s \{0, 1\}^n$
 $\text{ServerTokens}[id_{tk}] \leftarrow tk$
 $\hat{k} \leftarrow \{0, 1\}^\lambda$
 $etk \leftarrow \text{Enc}_1(\hat{k}, 0^{|tk|}, id_{tk})$

$tk \leftarrow \text{Dec}_1(\hat{k}, etk, id_{tk})$
if $|m| = \ell(n)$ **and** $P(m) = 1$
 $req \leftarrow \text{Enc}_2(tk, m_1, nc)$
else $req \leftarrow \perp$
 $C := C \cup \{req\}$
return (id_{tk}, etk, req)

OProcess(id_{tk}, req)

assert $\text{RespUsed}[id_{tk}] = 0$
assert $req \notin C$
 $tk \leftarrow \text{ServerTokens}[id_{tk}]$
 $m \leftarrow \perp$
 $\text{ServerTokens}[id_{tk}] \leftarrow \perp$
return m

□