

CS-E4340 Cryptography: Exercise Sheet 9

Estuardo Alpirez Bock, Christopher Brzuska and Kirthivaasan Puniamurthy

Aalto University, Finland

Abstract. This exercise deals with definitional studies for white-box cryptography and leakage resilient cryptography. In Definition 1, we provide a general definition for a white-box compiler. This compiler takes as input the secret key of a symmetric encryption scheme, and outputs a white-box encryption program. The white-box encryption program takes as input a message m and a nonce nc , and outputs an encryption of m . The white-box program is functionally equivalent to $\text{enc}(k, m, nc)$.

Exercise 1 (Side-channel analysis). Below we see a description of the Montgomery Ladder using projective coordinates. This algorithm is a popular choice for implementing the scalar multiplication in elliptic curve cryptosystems implemented in hardware. We recall that in such cryptosystems, the scalar k corresponds to the secret key and the multiplication kP is performed, e.g., for generating a signature or performing encryptions. The algorithm goes through an initialization phase (line 1), where 4 registers are set with some initial values. The algorithm then enters a loop (lines 2 to 10). The loop is executed once for each bit of the scalar k and depending on the bit value we execute either lines 4 and 5 or lines 7 and 8. Note that in both cases, the loop consists of the same arithmetic operations (multiplication, addition, squaring and overwriting the value of a register).

Figure 1 shows two power traces obtained during the execution of the algorithm (these power traces correspond actually to the initialization period and to the first 4 loop iterations of the algorithm). For the power trace above, a scalar $k = \kappa_1$ was used, while for the power trace below, a different scalar $k = \kappa_2$ was used. From observing these power traces, are you able to extract any bits from any of the two scalars?

Algorithm 1 Montgomery algorithm for the kP -operation using projective coordinates

Input: $k = (k_{l-1}, \dots, k_1, k_0)_2$ with $k_{l-1} = 1$, $P = (x, y) \in E(GF(2^m))$.

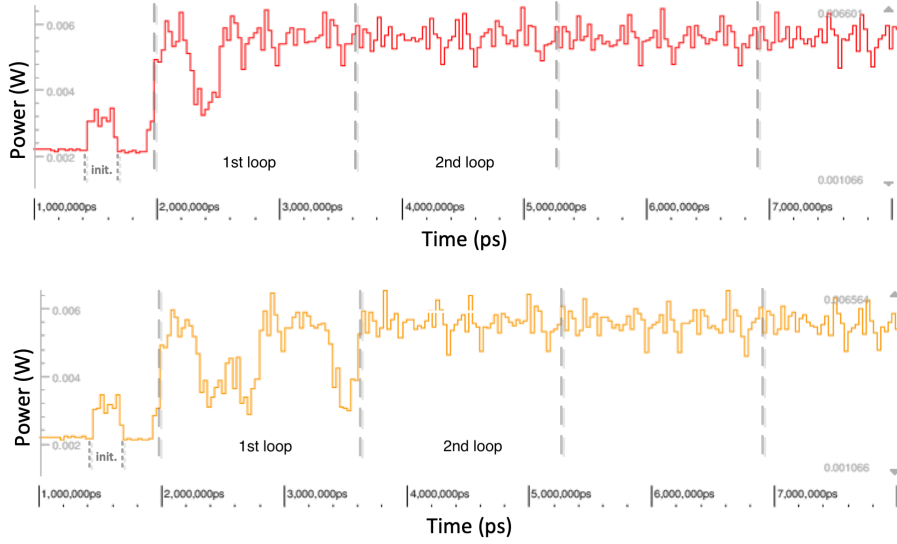
Output: $kP = (x_1, y_1)$.

```

1:  $X_1 \leftarrow x, Z_1 \leftarrow 1, X_2 \leftarrow x^4 + b, Z_2 \leftarrow x^2$ .
2: for  $i$  from  $l - 2$  downto 0 do
3:   if  $k_i = 1$  then
4:      $T \leftarrow Z_1, Z_1 \leftarrow (X_1 Z_2 + X_2 Z_1)^2, X_1 \leftarrow x Z_1 + X_1 X_2 T Z_2$ ,
5:      $T \leftarrow X_2, X_2 \leftarrow X_2^4 + b Z_2^4, Z_2 \leftarrow T^2 Z_2^2$ .
6:   else
7:      $T \leftarrow Z_2, Z_2 \leftarrow (X_2 Z_1 + X_1 Z_2)^2, X_2 \leftarrow x Z_2 + X_1 X_2 T Z_1$ ,
8:      $T \leftarrow X_1, X_1 \leftarrow X_1^4 + b Z_1^4, Z_1 \leftarrow T^2 Z_1^2$ .
9:   end if
10: end for
11:  $x_1 \leftarrow X_1 / Z_1$ .
12:  $y_1 \leftarrow y + (x + x_1)[X_1 + x Z_1](X_2 + x Z_2) + (x^2 + y)(Z_1 Z_2)] / (x Z_1 Z_2)$ .
13: return  $((x_1, y_1))$ .
```

Hints:

1. l is the length of the key, e.g. 233 bits. In this context, we say that the first bit of the key (k_{l-1}) is always a 1. The loop iterations start when we read the second bit (k_{l-2}).



2. $P = (x, y)$ is a point on the elliptic curve, the value of its coordinate x is used in line 1 to initialize some registers (e.g. register X_1). b is a coefficient used for defining the elliptic curve.
3. During the first loop iteration, the registers (or variables) X_1, X_2, Z_1 and Z_2 will be initialized with the values set on line 1.
4. Note that the Montgomery Ladder is expected to provide robustness against simple power analysis attacks. Thus, you are not really expected to extract *many* bits just from observing these two power traces. Note that most loop iterations look very similar and it is difficult to determine any key dependencies. However the first loop iteration looks different than the rest. Moreover, the first loop iteration for both power traces looks very different.
5. Recall that the power consumption of a device depends on the number of gates being switched during the computation. Thus, power consumption (per clock cycle) is highly dependent on the type of operations we are performing and also on the inputs we are processing. Do you think all multiplications demand equally high power? Or does it depend on the values we are multiplying?

Solution 1. For both power traces, we are only able to extract the key bit processed during the first loop. For the upper power trace, a key bit with value 1 was processed. For the lower trace, a key bit with value 0 was processed.

The reason we are able to distinguish is because of how the registers are initialized on line 1 of the algorithm. As we can see, the register Z_1 is initialized with a value 1. In both possible loop iterations, Z_1 is used for performing one or more multiplications. However any multiplication performed with register $Z_1 = 1$ will not be an expensive operation (i.e. it is a multiplication times 1), contrary to a multiplication performed with two values bigger than 1.

In the loop iteration for a key bit with value 1, the register Z_1 is multiplied once with the register X_2 , then Z_1 is overwritten. That is, we only have one multiplication with $Z_1 = 1$ for this loop iteration, hence the one dip in the upper power trace.

In the loop iteration for a key bit with value 0, the register Z_1 is multiplied once with the register X_2 , then it is multiplied with $X_1 X_2 T$, it is squared twice, multiplied with b and only in the end the register is overwritten. Thus we have more than one multiplication

performed with a value equal to 1, and we will have more than one operation which consumes a notably low amount of power. Hence we obtain more of such dips

Definition (White-box compiler). An algorithm **Comp** is called a white-box compiler for a symmetric encryption scheme $\text{se} = (\text{se.enc}, \text{se.dec})$ if for all key values $k \in \{0, 1\}^\lambda$, all nonces $nc \in \{0, 1\}^\lambda$ and all messages $m \in \{0, 1\}^*$, we have $\Pr[\text{se.enc}(k, m, nc) = \text{WB}(m, nc)] = 1$, where the probability is taken over $\text{WB} \leftarrow \$ \text{Comp}(k)$.

Exercise 2 (Security against key extraction). In the lecture we discussed how white-box programs may be susceptible to key-extraction attacks. Achieving security against key extraction is crucial for a white-box program. However as we will see in this exercise, white-box programs need to achieve further properties in order to provide any meaningful security.

Below we provide a definition capturing the property of security against key extraction for white-box programs.¹ Explain that while this definition captures security against key extraction, it does not really capture useful properties we would want from a white-box (cryptographic) program. Give an example of a program which achieves security against key extraction (i.e. is secure in the model below), yet it does not provide any further meaningful security (e.g. confidentiality of plaintexts). For this, you may design a secure symmetric-encryption scheme $\text{se}_{\text{example}}$ (assuming an existing secure symmetric-encryption scheme se) and a compiler **comp** for $\text{se}_{\text{example}}$ such that the compiler for $\text{se}_{\text{example}}$ outputs a white-box encryption program which is (1) functional equivalent to the encryption of $\text{se}_{\text{example}}$, (2) is secure against key extraction and (3) does not achieve any meaningful security in the presence of a white-box adversary.

Hint: Note that the definition below follows a notation style similar to the security definition for one-way functions (see Section 3.1 on the DAF (Security definition 1)).

Definition (Security against key extraction). Let $\text{Exp}_{\text{WB}, \mathcal{A}}^{\text{Kextr}}(1^\lambda)$ be a security experiment defined by

$$\begin{array}{l} \text{Exp}_{\text{Comp}, \mathcal{A}}^{\text{Kextr}}(1^\lambda) \\ k \leftarrow \$ \{0, 1\}^\lambda \\ \text{WB} \leftarrow \$ \text{Comp}(k) \\ k^* \leftarrow \$ \mathcal{A}(\text{WB}) \\ \text{if } k^* = k \text{ then} \\ \quad \text{return } 1 \\ \text{return } 0 \end{array}$$

A white-box compiler **Comp** for a symmetric encryption scheme se is secure against key extraction if for all PPT adversaries \mathcal{A} the winning probability $\text{Win}_{\text{Comp}, \mathcal{A}}^{\text{Kextr}}(1^\lambda) := \Pr[1 = \text{Exp}_{\text{Comp}, \mathcal{A}}^{\text{Kextr}}(1^\lambda)]$ is negligible in λ . (See discussion of search games in the DAF.)

Hints:

1. It may be useful to think of other primitives (besides encryption programs) which we may want to white-box. For instance, suppose you want to white-box a MAC program. As we recall from previous lectures, a good MAC should achieve an unforgeability property. Suppose we white-box a *bad* MAC program which does not achieve the unforgeability property. Suppose however, that the white-box version of this program manages to hide its key such

¹ In the literature, this definition might be referred to as *Unbreakability*. The name Unbreakability derives from the fact that when we manage to extract the key of a cryptographic implementation, we say that we *break* the implementation.

that it is indeed very difficult for an adversary to recover it from the implementation. Would such a MAC program be considered secure in the definition above, even though we know that it is a bad MAC?

Solution 2. We want to show that even if a white-box program hides its symmetric encryption key, an adversary might still be able to break further security properties we would usually desire. We simply define a symmetric encryption key that uses only half of its secret key. Then, we can throw away the other half (so that the key cannot be recovered since half of it is lost), but the scheme is still terribly insecure. It is impossible for the adversary to recover the complete key since only half of the key is located within the implementation and the other half is not really used.

Another example would be a scheme which only encrypts one half of the message and outputs the other half of the message in plain. Such a scheme would not achieve any notion of confidentiality. But if the adversary is unable to recover the secret key from the white-box program, the scheme would be considered secure in the definition above.

Exercise 3 (Leakage resilient cryptography). In wake of the evident threat of side-channel attacks, a research area known as *leakage resilient cryptography* was introduced. This area covers definitional studies considering an adversary who obtains some leakage of the secret key used within a cryptographic scheme. Then, the adversary tries to break security of the cryptographic scheme, given the leakage.

- a) Provide a definition for IND-CPA security of an encryption scheme under leakage. This definition should look similar to the traditional IND-CPA definition we have seen in previous lectures, with some additional steps where the adversary obtains leakage from the secret key.
- b) What do you think may be practical ways of achieving security under leakage?

Hints:

1. In definitions for cryptographic schemes secure under leakage, the adversary usually defines his own leakage function. However for the definition to make sense, the leakage function should have some restrictions. E.g. the leakage function should not be able to provide the adversary with the complete value of the secret key, else the adversary would always trivially win the security games. I.e. there may be some output length restrictions for the leakage functions. Alternatively, there might not be length restrictions for the output of the leakage functions, but the adversary should not be able to trivially derive the value of the secret key from the leakage. How would you capture these properties of the leakage function within the definition? What would be a reasonable output length restriction?

Solution 3. There are different ways of capturing IND-CPA security under leakage. Usually we have an adversary defining a leakage function in the beginning of the game, then the leakage function is applied to the secret key and its output is given to the adversary. The leakage function can be randomized, e.g. it can be a function which randomly selects some bits of the secret key and just returns them to the adversary. Then, given the leakage, the adversary proceeds to play a normal IND-CPA game. At some part of the game, we need to check whether the leakage is valid. I.e. whether the leakage is not too long or whether the secret key still has enough min-entropy, even given the leakage. Below we have an example of such a security definition. The adversary defines the leakage function and then the leakage function is applied to the key k . Then, the leakage is given

to the adversary and the adversary plays a classic IND-CPA game. In the third to last line, we check if the secret key still has enough min-entropy, given the leakage.

$\text{IND} - \text{CPA}_{\mu, \mathcal{A}}(1^n)$	$O(m_0, m_1)$
$b \leftarrow \$ \{0, 1\}$ $k \leftarrow \$ \text{kgen}(1^n)$ $\text{leak}(\cdot) \leftarrow \$ \mathcal{A}(1^n)$ $lk \leftarrow \$ \text{leak}(k)$ $b' \leftarrow \$ \mathcal{A}^O(lk)$ if $b' = b$ and $(\text{leak}(k) H_\infty(k)) \geq \mu$ return 1 else return 0	$c_b \leftarrow \$ \text{Enc}(k, m_b)$ return c_b

A trivial, practical way of increasing the robustness against leakage adversaries is by using keys of very large size. If the key is too large, the adversary might not be able to extract enough bits in order to break security. However the use of large keys within cryptographic schemes should not affect the efficiency of the schemes. One common strategy is to have a very large key but use only some randomly chosen bits from it for each encryption.

Exercise 4 (Impossibility). AES has been shown to be tremendously difficult to white-box. As we will see via an example in this exercise, some symmetric encryption schemes are actually impossible to be securely white-boxed, and they might not even achieve the modest goal of security against key extraction. Let se be an IND-CPA secure symmetric encryption scheme, let f be a pseudorandom function and let se' be the symmetric encryption scheme as defined below. Let comp be a (correct) white-box compiler for se' , i.e. $\text{WB} \leftarrow \$ \text{Comp}(k'_{\text{se}})$.

Give the code of an efficient attacker \mathcal{A} such that $\text{Win}_{\text{Comp}, \mathcal{A}}^{\text{Kextr}}(1^\lambda) := \Pr[1 = \text{Exp}_{\text{Comp}, \mathcal{A}}^{\text{Kextr}}(1^\lambda)]$ is non-negligible. You only need to argue about the success probability informally.

Hint: A black-box program means that an adversary is able to obtain input/output pairs but not know the code/internals of the program. What happens when an adversary instead interacts with a white-boxed program, i.e. it knows the program's code/internals, although it cannot make sense of it?

$\text{kgen}'(1^\lambda)$	$\text{enc}'(k_{\text{se}'}, x, \text{nc})$	$\text{dec}'(k_{\text{se}'}, c, \text{nc})$
$k_{\text{SE}} \leftarrow \$ \{0, 1\}^\lambda$ $k_x \leftarrow \$ \{0, 1\}^\lambda$ $k_{\text{nc}} \leftarrow \$ \{0, 1\}^\lambda$ $k_{\text{se}'} \leftarrow k_{\text{SE}} k_x k_{\text{nc}}$ return $k_{\text{se}'}$	$C \leftarrow \text{PARSE}(x)$ $k_{\text{SE}} k_x k_{\text{nc}} \leftarrow k_{\text{se}'}$ $d \leftarrow 1$ for i from 1 to λ $x_i \leftarrow \text{PRF}(k_x, \text{bin}_\lambda(i) \text{nc})$ $\text{nc}_i \leftarrow \text{PRF}(k_{\text{nc}}, \text{bin}_\lambda(i) \text{nc})$ if $C(x_i, \text{nc}_i) = 0$ $\text{enc}(k_{\text{SE}}, x_i, \text{nc}_i)$ $d \leftarrow d \wedge 1$ else $d \leftarrow d \wedge 0$ $c' \leftarrow \text{enc}(k_{\text{SE}}, x, \text{nc})$ if $d = 0$ $c \leftarrow 0 c'$ else $c \leftarrow 1 c' k_{\text{se}'}$ return c	$d \tilde{c} \leftarrow c$ $k_{\text{SE}} k_x k_{\text{nc}} \leftarrow k_{\text{se}'}$ if $d = 0$ $x \leftarrow \text{dec}(k_{\text{SE}}, \tilde{c}, \text{nc})$ else $c' k' \leftarrow \tilde{c}$ if $k' = k_{\text{se}'}$ $x \leftarrow \text{dec}(k_{\text{SE}}, c', \text{nc})$ else $x \leftarrow \perp$ return x

Solution 4. The key observation is that the encryption scheme outputs its secret key k , if the message given as input corresponds to a circuit describing the encryption algorithm (including its secret key). In the black-box attack model, an adversary cannot provide such an input since he does not have access to such a circuit, the adversary only sees the input/output behaviour of the program. Thus, this scheme is secure in the black-box attack model. In the white-box attack model however, the adversary gets access to the code describing the encryption program (even if he does not understand what's written there). Thus the adversary can copy that code, parse it into a message and provide it as input to the encryption program. The encryption program will leak its symmetric key. The adversary proceeds as follows.

```

A(WB)


---


 $x \leftarrow \text{PARSE}(\text{WB})$  // use the code describing WB as the message
 $\text{nc} \leftarrow \$_{ \{0,1\}^\lambda }$ 
 $c \leftarrow \text{WB}(x, \text{nc})$ 
 $1||c'|||k_{\text{se}'} \leftarrow c$ 
 $k_{\text{SE}}||k_x||k_{\text{nc}} \leftarrow k_{\text{se}'}$ 
return  $k_{\text{SE}}$ 

```

It follows from the correctness of se' that the adversary wins the key extraction game with overwhelming probability.

Exercise 5 (White-box as PKE). One possible view on security of white-box cryptography (without hardware-binding) is to build public key encryption scheme from a white-box of symmetric encryption scheme and demand that the resulting PKE satisfies IND-CPA security of a PKE. Given a white-box symmetric encryption scheme, describe a pke based on it and argue intuitively why it satisfies IND-CPA of PKE.

Solution 5. Key generation samples a random key, runs the compiler, outputs the key k as secret key and the white-box program as public-key. Encryption uses the white-box program. Decryption just uses k . Anybody can encrypt using the publicly known white-box program, but they cannot decrypt since they cannot extract the key k from it. Using IND-CPA security as a security definition of the compiler now demands that the white-box encryption program does not harm confidentiality of ciphertexts and, in particular, that the attacker cannot use the white-box encryption program for decrypting. This clearly implies that an adversary should not be able to extract the secret key from the white-box program (else, he could use that key for decrypting ciphertexts and breaking confidentiality).

Exercise 6 (Hardware-binding). In the lecture we discussed the issue that white-box programs might be susceptible to re-distribution attacks (or *code-lifting* attacks). In such attacks, an adversary copies the complete white-box program and simply uses it on a device of its choice. One method to mitigate such attacks is to compile the white-box program such that it can only be executed on one specific device. We refer to this property as hardware-binding. One way to implement hardware-binding is by having a program which tries to check on which device it is running before performing a computation. For instance, if the device running the white-box has a secure hardware module (which is not accessible to the white-box adversary), we could use one functionality of that module to check if we are running on the desired hardware.

Consider now a key derivation function (KDF) that we wish to white-box. Here, we consider a simple key derivation function which uses a secret key k_{kdf} to derive further keys from a context value e . The derived keys should be indistinguishable from random and thus, the syntax and security definition of this KDF are equal to those we have defined for PRFs². In this exercise we will study a security definition for a white-box KDF with hardware-binding. Below, we first introduce the syntax for our hardware module. Then we introduce the syntax for the white-box KDF. Finally, we provide a definition for the security of a white-box KDF.

We observe that the white-box KDF works in collaboration from a response coming from the hardware module. Without this response, the white-box KDF cannot compute anything. The security definition has some subtle flaws. Namely, currently, an adversary can win this game. Can you spot the flaws, explain them and correct them?

Definition (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} is a PPT algorithm, and the algorithms $\text{SubKgen}_{\text{HW}}$, Resp_{HW} and Check_{SW} are deterministic algorithms with the following syntax:

$$\begin{aligned} k_{\text{HWM}} &\leftarrow \$ \text{kgen}_{\text{HW}}(1^n) & k_{\text{HWS}} &\leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWM}}, \text{label}) \\ \sigma &\leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWM}}, \text{label}, x) & b &\leftarrow \text{Check}_{\text{SW}}(k_{\text{HWS}}, x, \sigma), \end{aligned}$$

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

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

where the probability is over the sampling of $k_{\text{HWM}} \leftarrow \$ \text{kgen}_{\text{HW}}(1^n)$.

The purposes of the algorithms are:

- The algorithm kgen_{HW} is used within the HWM to generate a main key k_{HWM} , which is located in the HWM.
- The algorithm $\text{SubKgen}_{\text{HW}}$ is used within the HWM to generate a subkey k_{HWS} , which is used when compiling the whitebox program (the whitebox program can later use the key k_{HWS} to verify that it is running in the correct hardware.)
- The algorithm Resp_{HW} is run in the HWM. The whitebox program can ask the HWM to give a response on which hardware it is running, that invokes Resp_{HW} algorithm which gives the answer.
- The algorithm Check_{SW} is used within the whitebox program to check that the certificate given by Resp_{HW} was valid. To do this, Check_{SW} uses the subkey k_{HWS} .

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

$$\text{WKDF} \leftarrow \$ \text{Comp}(k_{\text{kdf}}, k_{\text{HWS}}).$$

For all genuine k_{HWM} , for all k_{kdf} , for all label , for all e , for all $k_{\text{HWS}} = \text{SubKgen}_{\text{HW}}(k_{\text{HWM}}, \text{label})$ and $\sigma = \text{Resp}_{\text{HW}}(k_{\text{HWM}}, \text{label}, e)$, we have

$$\Pr[\text{KDF}(k_{\text{kdf}}, e) = \text{WKDF}(e, \sigma)] = 1$$

where the probability is taken over compiling $\text{WKDF} \leftarrow \$ \text{Comp}(k_{\text{kdf}}, k_{\text{HWS}})$.

Definition. A hardware-bound white-box key derivation scheme WBHW is IND-WKDF-secure if the real and ideal games $\text{Gind-wkdf}_{\text{WBHW}}^b$ are computationally indistinguishable, that is, for all PPT adversaries \mathcal{A} the advantage

$$\text{Adv}_{\mathcal{A}}^{\text{Gind-wkdf}_{\text{WBHW}}^0, \text{Gind-wkdf}_{\text{WBHW}}^1}(\lambda) := \left| \Pr[1 = \mathcal{A} \circ \text{Gind-wkdf}_{\text{WBHW}}^0] - \Pr[1 = \mathcal{A} \circ \text{Gind-wkdf}_{\text{WBHW}}^1] \right|$$

is negligible.

² We here use the term KDF, because it makes more sense conceptually, but if it is confusing, think PRF each time you read KDF.

Gind-wkdf⁰_{WBHW}

Parameters

λ : security parameter
WBHW: WKDF scheme

Package State

k_{kdf} : KDF key
 k_{HWm} : hardware main key
 k_{HWS} : hardware sub-key

GETWB(*label*)

assert WKDF $\neq \perp$
 $k_{\text{kdf}} \leftarrow_{\$} \{0, 1\}^\lambda$
 $k_{\text{HWm}} \leftarrow_{\$} \text{kgen}_{\text{HW}}(1^\lambda)$
 $k_{\text{HWS}} \leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWm}}, \text{label})$
WKDF $\leftarrow_{\$} \text{Comp}(k_{\text{kdf}}, k_{\text{HWS}})$
return WKDF

KDF()

$e \leftarrow_{\$} \{0, 1\}^\lambda$
 $\hat{k} \leftarrow \text{KDF}(k_{\text{kdf}}, e)$
return (\hat{k}, e)

Resp(*e*, *label*)

$\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWm}}, \text{label}, e)$
return σ

Gind-wkdf¹_{WBHW}

Parameters

λ : security parameter
WBHW: WKDF scheme

Package State

k_{kdf} : KDF key
 k_{HWm} : hardware main key
 k_{HWS} : hardware sub-key

GETWB(*label*)

assert WKDF $\neq \perp$
 $k_{\text{kdf}} \leftarrow_{\$} \{0, 1\}^\lambda$
 $k_{\text{HWm}} \leftarrow_{\$} \text{kgen}_{\text{HW}}(1^\lambda)$
 $k_{\text{HWS}} \leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWm}}, \text{label})$
WKDF $\leftarrow_{\$} \text{Comp}(k_{\text{kdf}}, k_{\text{HWS}})$
return WKDF

KDF()

$e \leftarrow_{\$} \{0, 1\}^\lambda$
 $\hat{k} \leftarrow_{\$} \{0, 1\}^\lambda$
return (\hat{k}, e)

Resp(*e*, *label*)

$\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWm}}, \text{label}, e)$
return σ

Solution 6. The whiteboxed KDF takes as input e and σ where σ is the response from the hardware module that tells whether the hardware is correct.

Now, if the adversary wants to win the game, the adversary should be able to produce a fake σ that convinces the whiteboxed KDF even though σ was generated by the adversary and not by the hardware module.

The current security definition allows the adversary to simply query the correct σ from Resp oracle.

```

 $\mathcal{A}$ 
-----
WKDF  $\leftarrow$  GETWB( $l$ )
( $k, e$ )  $\leftarrow$  KDF
 $\sigma \leftarrow$  Resp( $e, l$ )
if  $k = \text{WKDF}(e, \sigma)$ 
  return 1
else
  return 0

```

This adversary always returns 1 in the real game, while in the ideal game the key k is random, and hence adversary returns 1 with probability $2^{-\lambda}$, so the adversary has advantage $1 - 2^{-\lambda}$ which is non-negligible (almost one!).

We can easily modify the game by integrating a set which keeps track of the context values which have been generated within the KDF oracles. If the adversary tries to get a response with a context value which was generated within the oracle, then the oracle should return an error.

Gind-wkdf_{WBHW}⁰Parameters

λ : security parameter
 WBHW: WKDF scheme

Package State

k_{kdf} : KDF key
 k_{HWm} : hardware main key
 k_{HWs} : hardware sub-key
 Q : context set

GETWB(*label*)

assert $\text{WKDF} \neq \perp$
 $k_{\text{kdf}} \leftarrow \$ \{0, 1\}^\lambda$
 $k_{\text{HWm}} \leftarrow \$ \text{kgen}_{\text{HW}}(1^\lambda)$
 $k_{\text{HWs}} \leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWm}}, \text{label})$
 $\text{WKDF} \leftarrow \$ \text{Comp}(k_{\text{kdf}}, k_{\text{HWs}})$
return WKDF

KDF()

$e \leftarrow \$ \{0, 1\}^\lambda$
 $Q \leftarrow Q \cup \{e\}$
 $\hat{k} \leftarrow \text{KDF}(k_{\text{kdf}}, e)$
return (\hat{k}, e)

Resp(*e, label*)

assert $e \notin Q$
 $Q \leftarrow Q \cup \{e\}$
 $\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWm}}, \text{label}, e)$
return σ

Gind-wkdf_{WBHW}¹Parameters

λ : security parameter
 WBHW: WKDF scheme

Package State

k_{kdf} : KDF key
 k_{HWm} : hardware main key
 k_{HWs} : hardware sub-key
 Q : context set

GETWB(*label*)

assert $\text{WKDF} \neq \perp$
 $k_{\text{kdf}} \leftarrow \$ \{0, 1\}^\lambda$
 $k_{\text{HWm}} \leftarrow \$ \text{kgen}_{\text{HW}}(1^\lambda)$
 $k_{\text{HWs}} \leftarrow \text{SubKgen}_{\text{HW}}(k_{\text{HWm}}, \text{label})$
 $\text{WKDF} \leftarrow \$ \text{Comp}(k_{\text{kdf}}, k_{\text{HWs}})$
return WKDF

KDF()

$e \leftarrow \$ \{0, 1\}^\lambda$
 $Q \leftarrow Q \cup \{e\}$
 $\hat{k} \leftarrow \$ \{0, 1\}^\lambda$
return (\hat{k}, e)

Resp(*e, label*)

assert $e \notin Q$
 $Q \leftarrow Q \cup \{e\}$
 $\sigma \leftarrow \text{Resp}_{\text{HW}}(k_{\text{HWm}}, \text{label}, e)$
return σ