

# Watermarking Cryptographic Programs Against Arbitrary Removal Strategies

Ryo Nishimaki\*

Daniel Wichs†

April 20, 2015

## Abstract

A watermarking scheme for programs embeds some information called a mark into a program while preserving its functionality. No adversary can remove the mark without damaging the functionality of the program. In this work, we study the problem of watermarking various cryptographic programs such as pseudorandom function (PRF) evaluation, decryption, and signing. For example, given a PRF key  $K$ , we create a marked program  $\tilde{C}$  that evaluates the PRF  $F(K, \cdot)$ . An adversary that gets  $\tilde{C}$  cannot come up with *any* program  $C^*$  in which the mark is removed but which still evaluates the PRF correctly on even a small fraction of the inputs.

The work of Barak et al. (CRYPTO'01 and J.ACM, 59(2)) shows that, assuming indistinguishability obfuscation (iO), such watermarking is *impossible* if the marked program  $\tilde{C}$  evaluates the original program with perfect correctness. In this work we show that, assuming iO, such watermarking is *possible* if the marked program  $\tilde{C}$  is allowed to err with even a negligible probability, which would be undetectable to the user. We construct such a watermarking scheme with a secret-marking key used to embed marks in programs, and a public-detection key that allows anyone to detect marks in programs. For our security definition, we assume that the adversary can get oracle access to the marking functionality.

We emphasize that our security notion of watermark non-removability considers arbitrary adversarial strategies to modify the marked program – for example, an adversary could obfuscate the marked program and this should not remove the mark. This is in contrast to the prior works, such as that of Nishimaki (EUROCRYPT '13), which only consider restricted removal strategies that preserve the original structure of the marked program (e.g., as a vector of group elements), but do not provide security against arbitrary strategies.

## 1 Introduction

Digital watermarking enables us to embed some special information called a *mark* into digital objects such as images, movies, music files, or software. We often call such objects *marked*. There are two basic requirements for watermarking. The first is that a marked object should not be significantly different from the original object. The second is that malicious entities should not be able to remove embedded marks without somehow “destroying” the object (e.g., modify an image beyond recognition).

There are many works on watermarking perceptual objects such as images, movies, music files, etc. Most of them do not give a rigorous theoretical treatment and their constructions are heuristic and ad-hoc. The work of [HMW07] proposes general and rigorous definitions for such schemes, and gives an in-depth exploration of the connections between such definitions, but does not provide any actual constructions.

---

\*NTT Secure Platform Laboratories, nishimaki.ryo@lab.ntt.co.jp. This work was done while the author is visiting Northeastern University.

†Northeastern University, wichs@ccs.neu.edu. Research supported by NSF grants CNS-1347350, CNS-1314722, CNS-1413964.

**Watermarking Cryptographic Programs.** In this work, we focus on watermarking *programs*, and specifically ones that are cryptographic in nature, such as programs evaluating a pseudorandom function (PRF) or implementing a signing or a decryption procedure. For concreteness, let us focus on the case of PRFs.

We define a watermarking scheme for a PRF family  $\{F(K, \cdot)\}$  to consist of procedures  $\text{Mark}(mk, \cdot)$  and  $\text{Detect}(dk, \cdot)$  with a secret-marking key  $mk$  and a public-detection key  $dk$ . Given a PRF key  $K$ , the marking procedure  $\tilde{C} \leftarrow \text{Mark}(mk, K)$  creates a marked PRF circuit  $\tilde{C}$  that evaluates  $F(K, \cdot)$ . Although we will see that we cannot achieve *perfect correctness* where  $\tilde{C}(x) = F(K, x)$  for all inputs  $x$ , we will be able to achieve *statistical correctness* where we allow a negligible error probability. The detection procedure  $\text{Detect}(dk, C')$  either outputs that a circuit is marked or unmarked.

For security, we consider a game where we choose a random PRF key  $K$  and give the adversary the marked circuit  $\tilde{C} \leftarrow \text{Mark}(mk, K)$ . Intuitively, we require that the adversary cannot come up with *any* program that correctly evaluates the PRF and removes the mark. More precisely, the adversary produces a circuit  $C^*$  and we insist that either: (1) the circuit  $C^*$  does not even approximately compute  $F(K, \cdot)$ , meaning that  $C^*(x) = F(K, x)$  on at most an  $\varepsilon$ -fraction of the inputs  $x$ , or (2) the detection procedure correctly detects that the circuit is marked by outputting  $\text{marked} \leftarrow \text{Detect}(dk, C^*)$ . The parameter  $\varepsilon$ , called the “approximation factor” and we can set it to some small constant such as 1% or even to any 1/poly fraction. During the attack, the adversary is also given the public-detection key  $dk$  and access to the marking oracle  $\text{Mark}(mk, \cdot)$  that he can query on arbitrary PRF keys of his choice.

**Variants.** We call the above type of watermarking “mark-embedding” to denote that it only distinguishes between marked and unmarked program. We also consider “message-embedding” watermarking where the marking procedure can be used to embed an arbitrary message into the program and the detection procedure should recover the message.

We can also consider a relaxation of the above definition to the symmetric-key setting, where there is a single secret watermarking key  $wk$  used by the marking and detection procedures, and weaker security definitions depending on whether the adversary has access to the marking and/or detection oracles (a detection oracle  $\text{Detect}(wk, \cdot)$  allows an adversary to test whether a circuit is marked or unmarked – in the public-key setting this is subsumed by giving the adversary the public-detection key  $dk$ ).

**Motivation.** To highlight the usefulness of watermarking cryptographic functions, we describe two applications as examples. However, we emphasize that the concept should have broader applicability beyond these examples.

Consider an automobile manufacturer that wants to put electronic locks on its cars; the car contains a PRF key  $K$  and can only be opened by running an identification protocol where it chooses a random input  $x$  and the user must respond with  $F(K, x)$ . When a car is sold to a new owner, the owner is given a software key (e.g., a smart-phone application) consisting of marked program  $\tilde{C}$  that evaluates the PRF  $F(K, \cdot)$  and is used to open the car. The mark can embed some identifying information such as the owner’s name and address. Even if the software key is stolen, the thief cannot create some new software that would open the car while removing information about the original owner.

As another example, we can create a traitor-tracing scheme that preserves anonymity without having a central registry of users. In this scenario there is a content provider that streams encrypted copyrighted content to users. When users register, they provide some identifying information and they get a marked decryption key with this information embedded, but the identifying information does not need to be stored in any central registry in the long term.<sup>1</sup> The content provider only stores an anonymous encryption key for

---

<sup>1</sup>Alternately, we can implement this using a two-party computation protocol, where the user has her identifying information as an input (perhaps signed by a third party) and gets the marked key as an output without the content provider learning anything about the user at any point in time.

each user. If any user posts a valid decryption key in public, the identifying information about that user can be recovered from the key.<sup>2</sup>

**Impossibility?** The work of Barak et al. [BGI<sup>+</sup>01, BGI<sup>+</sup>12] initiated the first theoretical study of program watermarking. They propose a game-based definition which is similar to the weakest definition we consider in this work (in the symmetric-key setting with no marking/detection oracles given to the adversary) but require perfect correctness. Unfortunately, they show that this definition is unachievable assuming the existence of indistinguishability obfuscation.

The main intuition behind the negative result is to consider an attacker that takes a marked program and applies indistinguishability obfuscation (iO) to it. If the marked program implements the original program with *perfect correctness* then the result of applying iO to it should be indistinguishable from that of applying iO to the original program. Since the latter is unlikely to be marked, the same should apply to the former. Therefore, this presents a valid attack against watermarking in general.

Barak et al. note that the above attack crucially relies on the *perfect* (rather than merely *statistical*) correctness of the marked program, meaning that it correctly evaluates the original program on every input. They mention that otherwise “it seems that obfuscators would be useful in constructing watermarking schemes, because a watermark could be embedded by changing the value of the function at a random input, after which an obfuscator is used to hide this change.” This idea was not explored further in [BGI<sup>+</sup>01, BGI<sup>+</sup>12] and it is far from clear if a restricted notion of obfuscation such as iO would be sufficient and what type of watermarking/security can be achieved with this approach. Indeed, this idea serves as the starting point of our work.

**Other Related Work.** We mention that there are several other works [NSS99, YF11, Nis13] that propose concrete schemes for watermarking cryptographic functions, under several different definitions and assumptions. For example, the work of Nishimaki [Nis13] gives formal definitions and provably secure constructions for watermarking cryptographic functions (such as trapdoor functions). The main aspect that sets our work apart from these works is that they only consider restricted attacks which attempt to remove a watermark by outputting a new program which has some specific format (rather than an arbitrary program). For example, in the work of Nishimaki [Nis13], cryptographic keys (both marked and unmarked ones) are represented as a vector of group elements, and a valid attack is required to output a key of the same type. In particular, for all of the schemes proposed in these works, the watermark can be removed via the attack described in [BGI<sup>+</sup>01, BGI<sup>+</sup>12] where an adversary uses iO to obfuscate the marked program so as to preserve its functionality but completely change its structure.

## 1.1 Our Results and Techniques

We start by giving new formal definitions of program watermarking, along the lines of what we described earlier. To avoid the [BGI<sup>+</sup>01, BGI<sup>+</sup>12] impossibility results, our definition allows for statistical rather than perfect correctness.

As our main technical contribution, we show how to watermark a family of PRFs. Our scheme has a public-key detection procedure and achieves security in the presence of a marking oracle. We get a mark-embedding scheme (no message) that allows for any  $\varepsilon = 1/\text{poly}$  approximation factor and a message-embedding scheme that allows for approximation factors  $\varepsilon = 1/\sqrt{2} + 1/\text{poly}$ . In the case of message-

---

<sup>2</sup>Note that, in this scenario, we assume that the content is encrypted separately under each user’s key, which corresponds to a naive traitor-tracing scheme. Unfortunately, our watermarking definitions/constructions do not guarantee any security if the same key is marked with different messages and therefore each user would need to get a different key. The novelty here is that the identifying information about users is embedded in their keys rather than stored in a central registry. It would be interesting to combine this idea with more advanced traitor-tracing schemes where content doesn’t need to be encrypted separately.

embedding constructions, we show that there is an inherent lower bound of  $\varepsilon = 1/2$ . Both schemes rely on (sub-exponentially secure) indistinguishability obfuscation (iO).

For intuition, we also first construct a simple “basic scheme” which is a mark-embedding scheme (no message) in the symmetric-key setting (both the marking and detection keys are secret) and where the adversary does not get access to the marking or detection oracles. This scheme highlights the high-level ideas but removes much of the complexity needed for our full construction.

Once we have shown how to watermark a family of PRFs, we can easily extend this to watermarking other cryptographic primitives such as the decryption procedure of a public-key encryption or the signing procedure of a signature scheme. To do so, we rely on recent (obfuscation-based) constructions of public-key encryption and signatures where the decryption/signing procedures are simply PRF evaluation [SW14].

**Our techniques.** On a high level, to watermark a PRF  $F(K, \cdot)$ , we create a circuit  $C$  that evaluates the PRF correctly on almost all inputs  $x$ , except for some special set of “marked-points” (of negligible density) on which the program outputs incorrect values that can be detected. We obfuscate the circuit  $C$  to create a marked circuit  $\tilde{C}$  which “hides” the marked-points. The detection function can test an arbitrary circuit  $C'$  by evaluating it on some subset of the marked-points to determine if it is marked or unmarked.

For our basic scheme, we can implement the above idea by choosing a polynomial-sized set of  $\ell$  marked-points  $x_1, \dots, x_\ell$  and corresponding outputs  $y_1, \dots, y_\ell$  uniformly at random. These values form a watermarking key  $wk$  which is needed to both embed marked and detect marks. The detection function tests if on *at least one* of the  $\ell$  points  $x_i$  the program evaluates to  $y_i$ , and if so it determines that the program is “marked” else “unmarked”. We show that this implements “mark-embedding” watermarking (no message) with security in the symmetric-key setting against an attacker that doesn’t get access to the marking or detection oracles. We require that the PRF is a *puncturable PRF* [BW13, BGI14, KPTZ13, SW14] and the obfuscation scheme satisfies iO security. The proof follows in two steps. Firstly, using a careful series of hybrids, we show that an adversary that gets the marked program cannot differentiate between the marked-points  $x_i$  and uniformly random points in the domain. This relies on the puncturing proof technique by Sahai and Waters [SW14]. Secondly, we show that this implies that the probability of an adversary being able to come up with a program which is correct on even an  $\varepsilon$ -fraction of inputs, but is not detected (differs from the marked program on all of the  $\ell$  marked-points) is negligible as long as  $\ell = O(\lambda/\varepsilon)$  where  $\lambda$  is the security parameter. Therefore, we can set the approximation factor to any  $\varepsilon = 1/\text{poly}(\lambda)$  by choosing a correspondingly larger  $\ell$ .

Unfortunately, the above scheme already becomes insecure if the adversary has access to a marking oracle *or* a detection oracle, let alone if we make the key  $wk$  public. Fixing this in our full scheme is a significant challenge and requires us to introduce a more complex variant of the above ideas where: (1) The set of marked-points is different for each marked circuit and depends on the key  $K$  being marked, while still allowing the detection procedure to find marked-points without knowing  $K$ , (2) The number of marked-points is super-polynomial and the detection procedure tests the circuit on a small random subset of  $\ell$  of them; furthermore we obfuscate the main components of the detection procedure and release them as a public-detection key. The proof of security of the full scheme relies on a delicate sequence of super-polynomially many hybrids, which requires us to rely on complexity leveraging, and therefore assume the sub-exponential security of iO and puncturable PRFs.

To embed a message in the marked program, we follow the above approach but ensure that the outputs of the marked circuit on the marked-points encode some information about the message which can be recovered by the detection procedure.

## 2 Preliminaries

### 2.1 Notations

For any  $n \in \mathbb{N} \setminus \{0\}$ , let  $[n]$  be the set  $\{1, \dots, n\}$ . A bold face lower-case letter denotes a vector like  $\mathbf{x} = (x_1, \dots, x_n)$ . For two strings  $x_1$  and  $x_2$ ,  $x_1 || x_2$  denotes a concatenation of  $x_1$  and  $x_2$ . For program (or circuit)  $C$ ,  $C[a, b, c, \dots]$  denotes that  $C$  contains the values  $a, b, c, \dots$  “hardwired” in its description. When  $D$  is a random variable or distribution,  $y \leftarrow D$  denote that  $y$  is randomly selected from  $D$  according to its distribution. If  $S$  is a set, then  $x \leftarrow S$  denotes that  $x$  is uniformly selected from  $S$ . The expression  $y := z$  denotes that  $y$  is set, defined or substituted by  $z$ . We say that function  $f : \mathbb{N} \rightarrow \mathbb{R}$  is negligible in  $\lambda \in \mathbb{N}$  if  $f(\lambda) = \lambda^{-\omega(1)}$ . Hereafter, we use  $f \leq \text{negl}(\lambda)$  to mean that  $f$  is negligible in  $\lambda$ .

Let  $\mathcal{X} = \{X_\lambda\}_{\lambda \in \mathbb{N}}$  and  $\mathcal{Y} = \{Y_\lambda\}_{\lambda \in \mathbb{N}}$  denote two ensembles of random variables indexed by  $\lambda \in \mathbb{N}$ . We define a distinguish probability between  $\mathcal{X}$  and  $\mathcal{Y}$  to be  $\mu(\lambda)$  if for all PPT algorithm  $D$ ,

$$|\Pr[D(X_\lambda) = 1] - \Pr[D(Y_\lambda) = 1]| \leq \mu(\lambda)$$

We write  $\mathcal{X} \stackrel{\text{c}}{\approx}_\mu \mathcal{Y}$  to denote this. If  $\mu$  is negligible, we just write  $\mathcal{X} \stackrel{\text{c}}{\approx} \mathcal{Y}$ .

### 2.2 Definitions

In this section, we review basic notions and definitions used in this paper.

**Obfuscation.** The notion of indistinguishability obfuscation (iO) was proposed by Barak et. al. [BGI<sup>+</sup>01, BGI<sup>+</sup>12] and the first candidate construction was proposed by Garg, Gentry, Halevi, Raykova, Sahai, and Waters [GGH<sup>+</sup>13].

**Definition 2.1 (Indistinguishability Obfuscation [BGI<sup>+</sup>12, GGH<sup>+</sup>13])** A PPT algorithm  $i\mathcal{O}$  is an indistinguishability obfuscator (iO) if it satisfies the following two conditions.

**Functionality:** For all security parameter  $\lambda \in \mathbb{N}$ , for all circuit  $C$  for all input  $x$ , it holds that

$$\Pr[C'(x) = C(x) \mid C' \leftarrow i\mathcal{O}(1^\lambda, C)] = 1.$$

**Indistinguishability:** For all PPT distinguisher  $\mathcal{D}$  and all circuit ensembles  $C_0 = \{C_\lambda^{(0)}\}_{\lambda \in \mathbb{N}}$  and  $C_1 = \{C_\lambda^{(1)}\}_{\lambda \in \mathbb{N}}$  such that  $\forall \lambda, x : C_\lambda^{(0)}(x) = C_\lambda^{(1)}(x)$  and  $|C_\lambda^{(0)}| = |C_\lambda^{(1)}|$  we have:

$$\left| \Pr[\mathcal{D}(i\mathcal{O}(1^\lambda, C_\lambda^{(0)})) = 1] - \Pr[\mathcal{D}(i\mathcal{O}(1^\lambda, C_\lambda^{(1)})) = 1] \right| \leq \mu(\lambda)$$

where  $\mu(\lambda)$  is negligible. If the distinguishing probability  $\mu(\lambda)$  is bounded by  $2^{-\lambda^\epsilon}$  for some constant  $\epsilon$ , then we say  $i\mathcal{O}$  is sub-exponentially secure iO.

For simplicity, we write  $i\mathcal{O}(C)$  instead of  $i\mathcal{O}(1^\lambda, C)$  when the security parameter  $\lambda$  is clear from context.

**Pseudorandom Functions.** We review several variants of pseudorandom functions (PRFs).

**Definition 2.2 (Pseudorandom Functions)** A pseudorandom function  $\mathcal{F}$  consists of two PPT algorithms  $\mathcal{F} = (\text{Key}, F)$  and a pair of poly-time computable functions  $n(\cdot)$  and  $m(\cdot)$  that satisfy the following condition. For all PPT adversary  $\mathcal{A}$  and  $K \leftarrow \text{Key}(1^\lambda)$ , it holds

$$\left| \Pr[\mathcal{A}^{F(K, \cdot)} = 1] - \Pr[\mathcal{A}^{\mathcal{R}(\cdot)} = 1] \right| \leq \text{negl}(\lambda)$$

where  $F(K, \cdot) : \{0, 1\}^{n(\lambda)} \rightarrow \{0, 1\}^{m(\lambda)}$  is a deterministic function and  $\mathcal{R}$  is chosen uniformly at random from the set of all functions with the same domain/range.

The notion of puncturable pseudorandom function (pPRF) was proposed by Sahai and Waters [SW14].

**Definition 2.3 (Puncturable Pseudorandom Functions)** A puncturable pseudorandom function (pPRF)  $\mathcal{F}$  consists of three (probabilistic) algorithms  $\mathcal{F} = (\text{Key}, \text{Puncture}, F)$  and a pair of computable functions  $n(\cdot)$  and  $m(\cdot)$  that satisfy the following two conditions.

**Functionality preserving under puncturing:** For all polynomial size set  $S \subseteq \{0, 1\}^{n(\lambda)}$  and for all  $x \in \{0, 1\}^{n(\lambda)} \setminus S$ , it holds that

$$\Pr[F(K, x) = F(K \{S\}, x) \mid K \leftarrow \text{Key}(1^\lambda), K \{S\} := \text{Puncture}(K, S)] = 1.$$

**Pseudorandom at punctured points:** For all polynomial size set  $S = \{x_1, \dots, x_{k(\lambda)}\} \subseteq \{0, 1\}^{n(\lambda)}$  it holds that for all PPT adversary  $\mathcal{A}$ ,

$$\mu(\lambda) := \left| \Pr[\mathcal{A}(K \{S\}, \{F(K, x_i)\}_{i \in [k]}) = 1] - \Pr[\mathcal{A}(K \{S\}, U_{m(\lambda)|S}) = 1] \right| \leq \text{negl}(\lambda)$$

where  $K \leftarrow \text{Key}(1^\lambda)$ ,  $K \{S\} := \text{Puncture}(K, S)$  and  $U_\ell$  denotes the uniform distribution over  $\ell$  bits.

If the distinguish probability  $\mu(\lambda)$  is bounded by  $2^{-\lambda^\epsilon}$  for some constant  $\epsilon$ , then we say  $\mathcal{F}$  is sub-exponentially secure pPRF.

**Theorem 2.4 ([GGM86, BW13, BGI14, KPTZ13])** *If one-way functions exists, then for all efficiently computable  $n(\cdot)$  and  $m(\cdot)$ , there exists a pPRF family whose input is an  $n(\cdot)$  bit string and output is an  $m(\cdot)$  bit string.*

**Prefix Puncturable PRF.** We introduce the notion of prefix puncturable PRFs, which is an extension of pPRF. Prefix pPRFs allow us to puncture PRF keys at points that have a fixed prefix such as  $x^*$  where  $*$  means any (fixed length) string.

**Definition 2.5 (Prefix Puncturable Pseudorandom Functions)** A prefix pPRF  $\mathcal{F}$  consists of three (probabilistic) algorithms  $\mathcal{F} = (\text{Key}, \text{Puncture}, F)$  and poly-time computable functions  $n(\cdot)$ ,  $m(\cdot)$  that satisfy the following two conditions.

**Functionality preserving under prefix puncturing:** For all strings  $s$  and for all  $x \in \{0, 1\}^{n(\lambda)}$  such that  $x$  does not contains  $s$  as a prefix it holds that

$$\Pr[F(K, x) = F(K \{s\|*\}, x) \mid K \leftarrow \text{Key}(1^\lambda), K \{s\|*\} := \text{Puncture}(K, s\|*)] = 1.$$

**Pseudorandom at prefix punctured points:** For all PPT adversary  $\mathcal{A}$ , and for all values  $s$  and all polynomial sets of points  $S = \{x_1, \dots, x_{k(\lambda)}\} \subseteq \{0, 1\}^{n(\lambda)}$  such that each  $x_i \in S$  contains  $s$  as a prefix, it holds that

$$\mu(\lambda) := \left| \Pr[\mathcal{A}(K \{s\|*\}, \{F(K, x_i)\}_{i \in [k(\lambda)]}) = 1] - \Pr[\mathcal{A}(K \{s\|*\}, U_{m(\lambda)|S}) = 1] \right| \leq \text{negl}(\lambda)$$

where  $K \leftarrow \text{Key}(1^\lambda)$ ,  $K \{s\|*\} := \text{Puncture}(K, s\|*)$ , and  $U_\ell$  denotes the uniform distribution over  $\ell$  bits.

If the distinguish probability  $\mu(\lambda)$  is bounded by  $2^{-\lambda^\epsilon}$  for some constant  $\epsilon$ , then we say  $\mathcal{F}$  is sub-exponentially secure prefix pPRF.

We can easily construct prefix pPRFs from the GGM tree-based construction in the same way as puncturable PRFs.

**Injective PRFs.** We need pPRFs that satisfy the following injectivity condition.

**Definition 2.6 (Injective property of pPRF)**

$$\Pr_{K \leftarrow \text{Key}(1^\lambda)} [\exists \alpha, K' \neq K : F(K, \alpha) = F(K', \alpha)] \leq \text{negl}(\lambda)$$

In other words this says that w.o.p. over the choice of  $K$ , the value  $F(K, \alpha)$  uniquely determines  $K$ . The GGM tree-based pPRF from a special *injective* PRG satisfies this requirement. A special PRG is  $G(x) = G_0(x) \| G_1(x)$  where  $G_0(x)$  (resp,  $G_1(x)$ ) is the first (resp, last) half of the bits of  $G(x)$  and  $G_0$  and  $G_1$  are each *injective*, meaning that  $\Pr_x[\exists x' \neq x : G(x) = G(x')] \leq \text{negl}(\lambda)$ . Moreover each of  $|G_0(x)|, |G_1(x)|$  is of length at most  $|x| + p(\lambda)$  where  $p(\lambda)$  is some fixed polynomial. This ensures that, in the GGM tree-based pPRF construction, the size of the output can grow by some additive amount in each level but the overall growth in a level  $d$  tree is bounded by  $d \text{poly}(\lambda)$ . Such PRGs can be constructed from standard assumptions such as DDH or LWE.

**Injective pPRF from LWE.** For example, using the learning with errors (LWE) assumption, we define  $G_A : \mathbb{Z}_q^n \rightarrow \mathbb{Z}_p^m$  as  $G_A(x) := \lfloor A^T \cdot x \rfloor_p$  where operator  $\lfloor \cdot \rfloor_p$  returns the nearest integer (for each coordinate) modulo  $p$ . We can set  $q := p^2 = 2^{2k}$  for some  $k = O(\lambda)$  and  $m := 4n + O(\lambda)$ . Let  $A = A_0 \| A_1$  where  $A_0, A_1 \in \mathbb{Z}_q^{n \times m/2}$ , then  $G_b(x) = \lfloor A_b^T x \rfloor$ . In this case, each  $G_b(x)$  is injective w.o.p. over the choice of  $A$  and it maps  $2nk$  bits to  $2nk + O(k\lambda)$  bits. See [BPR12] for details about the LWE assumption and proof of security of the above construction.

**Injective pPRF from DDH.** Alternatively, it may seem that using DDH, we can set  $G_{g_1, g_2}(x) = g_1^x, g_2^x$  where  $g_1, g_2$  are generators of some group  $\mathbb{G}$  of primer order  $p$ . Unfortunately, the outputs cannot be directly used as PRG inputs in the next level of the tree since they are group elements rather than exponents and we do not know how to extract out two uniform values in  $\mathbb{Z}_p$  from them. Nevertheless, this approach can be made to work by defining  $G_{g_1, g_2, g_3, h_0, h_1}(x) = h_0(g_1^x, g_2^x, g_3^x), h_1(g_1^x, g_2^x, g_3^x)$  where  $h_0, h_1$  are universal hash functions that map  $\mathbb{G}^3 \rightarrow \mathbb{Z}_{p'}$  for some  $p'$  such that  $\log(p') = \log(p) + O(\lambda)$  and  $\log(p') \leq (3/2) \log(p) - \Omega(\lambda)$ . This ensures injectivity (we are hashing  $p$  balls into  $p'$  bins and therefore for any fixed ball there is unlikely to be another ball colliding with it). It also ensures pseudorandomness security by thinking of  $h_0, h_1$  as extractors via the leftover-hash lemma. In the context of the GGM construction we need a hierarchy of DDH groups of order  $p_1, p_2, \dots$  (one for each level) where  $\log(p_{i+1}) = \log(p_i) + O(\lambda)$ . Therefore the output does not get “too large”.

### 3 Definition of Watermarking

We begin by defining the notion of program watermarking. Our definition is similar to the game-based definition of Barak et al. [BGI<sup>+</sup>01, BGI<sup>+</sup>12] with the main difference that: (1) we allow “statistical” rather than perfect correctness, (2) the challenge circuit to be marked is chosen uniformly at random from the circuit family (for example, in the case of PRFs, this corresponds to marking a random PRF key), (2) we strengthen the definition to the public-key detection setting and give the attacker access to the marking oracle.

We start by defining two types of schemes which we will refer to as “mark-embedding” and “message-embedding” watermarking. In the former programs are simply either marked or unmarked while in the latter they can be marked with a (long) message. For conciseness, we define the notions simultaneously by defining the more general notion of “message-embedding” watermarking and thinking of “mark-embedding” watermarking as a special case where the message space  $\mathcal{M}$  consists of a single value  $\mathcal{M} = \{\text{marked}\}$ .

**Definition 3.1 (Watermarking Syntax)** A *message-embedding watermarking* scheme for a circuit class  $\{\mathcal{C}_\lambda\}_{\lambda \in \mathbb{N}}$  and a message space  $\mathcal{M} = \{\mathcal{M}_\lambda\}$  consists of three algorithms (Gen, Mark, Detect).

**Key Generation:**  $(dk, mk) \leftarrow \text{Gen}(1^\lambda)$  takes as input the security parameter and outputs a pair of keys  $(dk, mk)$ , called detection key and mark key.

**Mark:**  $\tilde{C} \leftarrow \text{Mark}(mk, C, \text{msg})$  takes as input a mark key, a circuit  $C \in \mathcal{C}_\lambda$  and a message  $\text{msg} \in \mathcal{M}_\lambda$  and outputs a marked circuit  $\tilde{C}$ .

**Detect:**  $\text{msg}' \leftarrow \text{Detect}(dk, C')$  takes as input a detection key and an arbitrary circuit  $C'$  (which may not be in the family  $\mathcal{C}_\lambda$ ) and outputs  $\text{msg}' \leftarrow \text{Detect}(dk, C')$  where  $\text{msg}' \in \mathcal{M} \cup \{\text{unmarked}\}$ . We say that detection is *black-box* if it only uses oracle-access to  $C'$  to evaluate it on various inputs, but ignores its description otherwise.

We also define *mark-embedding watermarking* which is the same as above but the message space  $\mathcal{M}$  consists of a single value  $\mathcal{M} = \{\text{marked}\}$ . In this case we omit the input  $\text{msg}$  from the syntax of the marking algorithm and simply write  $\tilde{C} \leftarrow \text{Mark}(mk, C)$ . We also define a variant with *symmetric-key detection* where there is a single secret watermarking key  $wk = mk = dk$ .

Before we define the correctness and security of watermarking, we define the following notion of *approximating* a function.

**Definition 3.2 ( $\varepsilon$ -Approximating a Function)** A circuit  $C$  is said to  $\varepsilon$ -approximate a function  $f : \{0, 1\}^n \rightarrow \{0, 1\}^*$ , denoted by  $C \cong_\varepsilon f$ , if  $\Pr_{x \leftarrow \{0, 1\}^n} [C(x) = f(x)] \geq \varepsilon$ .

We are now ready to define program watermarking.

**Definition 3.3 (Watermarking Security)** A watermarking scheme (Gen, Mark, Detect) for circuit family  $\{\mathcal{C}_\lambda\}_{\lambda \in \mathbb{N}}$  and with message space  $\mathcal{M} = \{\mathcal{M}_\lambda\}$  is required to satisfy the following properties.

**Statistical Correctness:** There is a negligible function  $\nu(\lambda)$  such that for any circuit  $C \in \mathcal{C}_\lambda$ , any message  $\text{msg} \in \mathcal{M}_\lambda$  and any input  $x$  in the domain of  $C$ , it holds that

$$\Pr \left[ \tilde{C}(x) = C(x) : \begin{array}{l} (dk, mk) \leftarrow \text{Gen}(1^\lambda) \\ \tilde{C} \leftarrow \text{Mark}(mk, C, \text{msg}) \end{array} \right] \geq 1 - \nu(\lambda).$$

**Detection Correctness:** For every  $C \in \mathcal{C}_\lambda$ ,  $\text{msg} \in \mathcal{M}_\lambda$  and  $(dk, mk) \leftarrow \text{Gen}(1^\lambda)$ :

$$\text{Detect}(dk, \text{Mark}(mk, C, \text{msg})) = \text{msg}.$$

**Meaningfulness:** For every circuit  $C$  (not necessarily in  $\mathcal{C}_\lambda$ ), it holds that

$$\Pr_{(dk, mk) \leftarrow \text{Gen}(1^\lambda)} [\text{Detect}(dk, C) \neq \text{unmarked}] \leq \text{negl}(\lambda).$$

**Non-Removability:** For every PPT  $\mathcal{A}$  we have

$$\text{Adv}_{\mathcal{A}}^{\text{nrmv}}(\lambda, \varepsilon) := \Pr[\text{Exp}_{\mathcal{A}}^{\text{nrmv}}(\lambda, \varepsilon) = 1] \leq \text{negl}(\lambda)$$

where  $\varepsilon$  is a parameter of the scheme called the *approximation factor* and  $\text{Exp}_{\mathcal{A}}^{\text{nrmv}}(\lambda, \varepsilon)$  is the game defined next.

**Definition 3.4 (Non-Removability Security Game)** The game  $\text{Exp}_{\mathcal{A}}^{\text{nrmv}}(\lambda, \varepsilon)$  is defined as follows.



1. The challenger generates  $(dk, mk) \leftarrow \text{Gen}(1^\lambda)$  and gives  $dk$  to the adversary  $\mathcal{A}$ .
2. The adversary has oracle access to the mark oracle  $\mathcal{MO}$ . If  $\mathcal{MO}$  is queried a circuit  $C_i$  and message  $\text{msg}_i$ , then it answers  $\text{Mark}(mk, C_i, \text{msg}_i)$ .
3. At some point, the adversary sends a challenge message  $\text{msg} \in \mathcal{M}_\lambda$  to the challenger. The challenger samples a circuit  $C \leftarrow \mathcal{C}_\lambda$  uniformly at random and generates  $\tilde{C} \leftarrow \text{Mark}(mk, C, \text{msg})$ , and sends  $\tilde{C}$  to the adversary.
4. The adversary can make more queries to  $\mathcal{MO}$  and finally outputs a circuit  $C^*$ . If it holds  $C^* \cong_\varepsilon \tilde{C} \wedge \text{Detect}(dk, C^*) \neq \text{msg}$  then the experiment outputs 1, otherwise 0. <sup>3</sup>

We can also weaken the above definition in several ways that we mention below.

**Variant: Selective-Message Security and Mark-Embedding.** Firstly, in the case of message-embedding watermarking, we could consider a *selective-message security* variant where the adversary has to declare the message  $\text{msg}$  embedded in the challenge circuit at the very beginning of the game, before seeing the detection key or making queries to the marking oracle. In this case the security game consists of the following steps:

1. The adversary sends a challenge message  $\text{msg} \in \mathcal{M}_\lambda$  to the challenger.
2. The challenger generates  $(dk, mk) \leftarrow \text{Gen}(1^\lambda)$ , samples a circuit  $C \leftarrow \mathcal{C}_\lambda$  uniformly at random and generates  $\tilde{C} \leftarrow \text{Mark}(mk, C, \text{msg})$ . It gives  $dk, \tilde{C}$  to the adversary  $\mathcal{A}$ .
3. The adversary  $\mathcal{A}$  gets oracle access to the mark oracle  $\mathcal{MO}$  and at some point outputs a circuit  $C^*$ . If it holds  $C^* \cong_\varepsilon \tilde{C} \wedge \text{Detect}(dk, C^*) \neq \text{msg}$  then the experiment outputs 1, otherwise 0.

Note that selective-message security implies full security via a guessing argument – the reduction loses a factor of  $|\mathcal{M}_\lambda|$  in the advantage and therefore this relies on complexity leveraging when  $|\mathcal{M}_\lambda|$  is super-polynomial.

In the case of *mark-embedding* watermarking, there is no challenge message (we always have  $\text{msg} = \text{marked}$ ) and therefore the game becomes as above with step 1 omitted.

**Variant: Symmetric-Key Schemes and No-Oracle Security.** We can also consider a variant in the symmetric-key setting where there is a single watermarking key  $wk = mk = dk$  and the attacker is not given  $dk$  in step 1 of the security game. In this case we can consider several variants of the definition depending on whether the adversary has access to the marking and/or detection oracles. The simplest setting we consider is a mark-embedding scheme in the symmetric-key setting with no-oracle security, where the attacker does not get access to either the marking or the detection oracles. In this case, the security game just consists of the following two steps:

1. The challenger chooses  $wk \leftarrow \text{Gen}(1^\lambda)$ ,  $C \leftarrow \mathcal{C}_\lambda$  and gives  $\tilde{C} \leftarrow \text{Mark}(wk, C)$  to the adversary.
2. The adversary outputs a circuit  $C^*$ . If  $C^* \cong_\varepsilon \tilde{C} \wedge \text{Detect}(wk, C^*) \neq \text{marked}$  then the experiment outputs 1, otherwise 0.

---

<sup>3</sup>The definition would be equivalent if we had required  $C^* \cong_\varepsilon C$  instead of  $C^* \cong_\varepsilon \tilde{C}$ , up to a negligible difference in  $\varepsilon$ , since by statistical correctness we have  $C \cong_\delta \tilde{C}$  for some  $\delta = 1 - \text{negl}(\lambda)$ .

**Note on Statistical Correctness.** We mention that, by an averaging argument, the statistical correctness requirement implies that for *any* distribution  $\mathcal{D}$  over inputs  $x$ , with overwhelming probability over the choice of the marked circuit  $\tilde{C}$ , we have  $\Pr_{x \leftarrow \mathcal{D}}[\tilde{C}(x) = C(x)] \geq 1 - \text{negl}(\lambda)$ . Therefore, this requirement is more meaningful than simply insisting that  $\tilde{C} \cong_\delta C$  for some  $\delta = 1 - \text{negl}(\lambda)$ .

**Note on the Circuit Family.** We also note that the marking procedure expects to get a circuit  $C \in \mathcal{C}_\lambda$  from the designated family. For example, in the case of watermarking a PRF family  $\{F(K, \cdot)\}$ , we will assume that the family  $\mathcal{C}_\lambda$  consists of some “canonical representation” of the circuits computing the function  $F(K, \cdot)$ , and we can simply equate circuits  $C \in \mathcal{C}_\lambda$  with PRF keys  $K$ ; we will write  $\text{Mark}(mk, K, \text{msg})$  in place of  $\text{Mark}(mk, C, \text{msg})$  where the  $C$  is the canonical circuit computing  $F(K, \cdot)$ . In the non-removability security game, the fact that the challenge circuit  $C \leftarrow \mathcal{C}_\lambda$  is chosen randomly corresponds to marking a random PRF key. However, we emphasize that the detection procedure  $\text{Detect}(dk, C')$  must work on arbitrary circuits  $C'$ , which may not be from the designated family, and the adversary wins if he can create *any* circuit  $C'$  which is correct and on which detection fails.

**Impossibility of Message-Embedding for  $\varepsilon \leq \frac{1}{2}$ .** Unfortunately, message-embedding watermarking is impossible in a setting where the adversary can make a single non-adaptive query to the *marking oracle* (even without a detection oracle) when  $\varepsilon \leq \frac{1}{2}$ . This holds even if the message space is just  $\mathcal{M} = \{0, 1\}$  and even in the symmetric-key setting when the adversary has no access to a detection oracle. The adversary chooses a random bit  $b \leftarrow \{0, 1\}$  at the beginning of the game. It makes a single non-adaptive call to the marking oracle on a random circuit  $C \leftarrow \mathcal{C}_\lambda$  with the message bit  $1 - b$  and gets back some marked circuit  $\tilde{C}_{1-b}$ . It then chooses the challenge message  $b$  and gets back some challenge marked circuit  $\tilde{C}_b$ . It creates a circuit  $C^*$  such that, on inputs  $x$  that start with the bit  $b$ ,  $C^*(x) = \tilde{C}_b(x)$  and on inputs  $x$  that start with the bit  $1 - b$ ,  $C^*(x) = \tilde{C}_{1-b}(x)$ . This ensures that  $C^* \cong_\varepsilon \tilde{C}_b$  for  $\varepsilon = \frac{1}{2}$  meaning that  $C^*$  is an  $\varepsilon$ -approximation of the challenge marked circuit. But since the way that  $C^*$  is created is symmetric with respect to  $b$  and  $1 - b$ , we know that  $\Pr[\text{Detect}(dk, C^*) = b] \leq \frac{1}{2}$ .

## 4 Watermarking PRFs

In this section, we construct schemes for watermarking any puncturable PRF family. To build up intuition, we start with a simple “basic construction” which is a mark-embedding symmetric-key watermarking scheme with “no-oracle security” (the attacker does not get access to either the marking or the detection oracles). We then explain the challenges that we need to overcome to get full security in the public-key detection setting, and give our full construction which achieves this.

For all of the schemes, let  $\mathcal{F} = (\text{Key}, \text{Puncture}, F)$  be some puncturable PRF (pPRF) family where, for  $K \leftarrow \text{Key}(1^\lambda)$  we have  $F(K, \cdot) : \mathcal{D}_\lambda \rightarrow \mathcal{R}_\lambda$  with  $\mathcal{D}_\lambda = \{0, 1\}^{n(\lambda)}$ , and  $\mathcal{R}_\lambda = \{0, 1\}^{m(\lambda)}$  for some super-logarithmic  $n(\lambda), m(\lambda) = \omega(\log \lambda)$ . We often drop  $\lambda$  from  $\mathcal{D}_\lambda$  and  $\mathcal{R}_\lambda$ . We construct a watermarking scheme for *PRF evaluation* of  $\mathcal{F}$ . We identify the PRF evaluation circuits computing the function  $F(K, \cdot)$  with their key  $K$  and assume (without loss of generality) that the marking procedure just takes  $K$  as an input.

### 4.1 Basic Construction: Symmetric-Key Watermarking, No Oracles

In this section, we start by proposing the most basic of our watermarking schemes, which is a mark-embedding construction in the symmetric-key setting and achieves security against an adversary with no oracle access (weakest security). We call this “the basic construction” throughout the paper, and the intuition for it was given in the introduction.

**Construction.** Our watermarking scheme with “approximation factor”  $\varepsilon = \varepsilon(\lambda)$  is defined as follows.

**Gen( $1^\lambda$ ):** Set  $\ell := \lambda/\varepsilon$ . Choose an  $\mathbf{x} = (x_1, \dots, x_\ell) \leftarrow \mathcal{D}^\ell$  and a vector  $\mathbf{y} = (y_1, \dots, y_\ell) \leftarrow \mathcal{R}^\ell$  uniformly at random and set  $wk := (\mathbf{x}, \mathbf{y})$ . We can assume that  $x_1, \dots, x_\ell$  are distinct.

**Mark( $wk, K$ ):** Output  $i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])$ . The program  $C[K, \mathbf{x}, \mathbf{y}]$  is described in Figure 1.

**Detect( $wk, C'$ ):** If there exists  $i \in [\ell]$  such that  $y_i = C'(x_i)$ , then output marked, else unmarked.

**Constants:**  $\mathbf{x} = (x_1, \dots, x_\ell) \in \mathcal{D}^\ell$ ,  $\mathbf{y} = (y_1, \dots, y_\ell) \in \mathcal{R}^\ell$ , pPRF key  $K$ .

**Inputs:**  $x \in \mathcal{D}$

**Code:** If there exists  $i$  such that  $x = x_i$ , then output  $y_i$ . Else, output  $F(K, x)$ .

Figure 1: Program  $C[K, \mathbf{x}, \mathbf{y}]$

We call the value  $\mathbf{x} = (x_1, \dots, x_\ell)$  the *marked-points* since marked programs outputs special values from  $\mathbf{y}$  for these inputs. We can easily verify that the basic construction satisfies statistical correctness, detection correctness, and meaningfulness properties in Definition 3.3.

**Security Proofs for the Basic Construction.** We prove the security, non-removability, of the basic construction.

**Theorem 4.1** *The basic construction satisfies symmetric-key non-removable security with no oracles (where the adversary in the security game of non-removability has oracle access to neither the mark nor detect oracle) if  $\mathcal{F}$  is a pPRF and  $i\mathcal{O}$  is an indistinguishability obfuscator.*

We prove Theorem 4.1 by using the following two lemmas.

**Lemma 4.2** *If  $\mathcal{F}$  is a pPRF and  $i\mathcal{O}$  is an indistinguishability obfuscator, then the following holds. For  $K \leftarrow \text{Key}(1^\lambda)$  and uniformly random and independent  $\mathbf{x}, \mathbf{x}' \leftarrow \mathcal{D}^\ell$  and  $\mathbf{y} \leftarrow \mathcal{R}^\ell$ , it holds that*

$$(\mathbf{x}, i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])) \stackrel{\approx}{\sim} (\mathbf{x}', i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])).$$

**Lemma 4.3** *If a PPT adversary  $\mathcal{A}$  that has no oracle access breaks the non-removability security with probability  $\delta(\lambda)$ , then we can construct a PPT distinguisher  $\mathcal{D}$  that can distinguish  $(\mathbf{x}, i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}]))$  and  $(\mathbf{x}', i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}]))$  with probability  $\delta(\lambda) - \text{negl}(\lambda)$ .*

Lemma 4.2 and 4.3 immediately imply Theorem 4.1.

*Proof of Lemma 4.3.* Given an input  $(\mathbf{x}^* = (x_1^*, \dots, x_\ell^*), \tilde{C} = i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}]))$ , we construct a distinguisher  $\mathcal{D}$  that uses the adversary  $\mathcal{A}$  to distinguish between  $\mathbf{x}^* = \mathbf{x}$  and  $\mathbf{x}^* = \mathbf{x}'$  for uniformly random and independent  $\mathbf{x}'$ . The distinguisher  $\mathcal{D}$  gives  $\tilde{C}$  to the adversary  $\mathcal{A}$  as a challenge marked program. Eventually,  $\mathcal{A}$  outputs a program  $C^*$  and  $\mathcal{D}$  verifies whether there exists  $i$  such that  $C^*(x_i^*) = \tilde{C}(x_i^*)$  or not. Formally,  $\mathcal{D}$  proceeds as follows.

1.  $\mathcal{D}$  gets  $(\mathbf{x}^*, \tilde{C})$  and gives  $\tilde{C}$  to  $\mathcal{A}$ .
2.  $\mathcal{D}$  receives  $C^*$  from  $\mathcal{A}$  and runs  $\text{Test}(C^*, \tilde{C})$  described in Figure 2 to test if  $C^*$  approximates  $\tilde{C}$ . If  $\text{Test}$  outputs 0, then  $\mathcal{D}$  outputs 0. Otherwise, goes to the next step.
3. If for all  $i \in [\ell]$ ,  $\tilde{C}(x_i^*) \neq C^*(x_i^*)$ , then  $\mathcal{D}$  outputs 1, otherwise 0.

**Inputs:** Two programs  $C_0$  and  $C_1$ .

**Parameters:**  $0 < \varepsilon < 1$ .

Set  $\text{cnt} := 0$  and  $R := 8\lambda(1/\varepsilon)^2$ . For  $i = 1, \dots, R$ , do

1. Choose  $z_i \leftarrow \mathcal{D}$
2. If  $C_0(z_i) = C_1(z_i)$  then set  $\text{cnt} := \text{cnt} + 1$

Let  $\hat{\varepsilon} := \text{cnt}/R \in [0, 1]$  be the fraction of trials in which  $C_0(z_i) = C_1(z_i)$ .

If  $\hat{\varepsilon} < \frac{3}{4}\varepsilon$ , then output 0, else 1.

Figure 2: Test algorithm Test for approximation of programs

**Analysis of Test.** Let  $Z_i$  be a random variable such that  $Z_i = 1$  if and only if  $C_0(z_i) = C_1(z_i)$  for  $z_i \leftarrow \mathcal{D}$ . Note that  $Z_i$  are independent and identically distributed.

If  $C_0 \cong_\varepsilon C_1$  then  $\Pr[Z_i = 1] \geq \varepsilon$  and the Chernoff bound implies  $\Pr[\hat{\varepsilon} < \frac{3}{4}\varepsilon] \leq \exp(-\frac{1}{8}R\varepsilon^2) \leq \exp(-\lambda)$ . This means that the test outputs 1 with overwhelming probability.

If it is *not* the case that  $C_0 \cong_{\varepsilon/2} C_1$  then  $\Pr[Z_i = 1] < \varepsilon/2$  and the Chernoff bound implies  $\Pr[\hat{\varepsilon} \geq \frac{3}{4}\varepsilon] \leq \exp(-\frac{1}{8}R\varepsilon^2) \leq \exp(-\lambda)$ . This means that the test output 0 with overwhelming probability.

**Analysis of  $\mathcal{D}$ .** Next, we analyze how  $\mathcal{D}$  works.

If the given distribution is the real one ( $\mathbf{x}, i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])$ ) (i.e.,  $\mathbf{x}^* = \mathbf{x}$ ), then  $\mathcal{A}$  wins with the non-removability game with probability  $\delta(\lambda)$ , meaning that it outputs an un-marked program  $C^*$  that  $\varepsilon$ -approximates  $\tilde{C}$ . If this happens, then Test passes (outputs 1) with overwhelming probability  $1 - \text{negl}(\lambda)$  and, for all  $i \in [\ell]$ ,  $C^*(x_i^*) \neq \tilde{C}(x_i^*)$ . Therefore, in this case,  $\mathcal{D}$  outputs 1 with probability at least  $\delta(\lambda) - \text{negl}(\lambda)$ .

If the given distribution is random one ( $\mathbf{x}', i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])$ ) (i.e.,  $\mathbf{x}^* = \mathbf{x}'$ ), then  $\mathbf{x}$  and  $\mathbf{x}'$  are uniformly random and independently distributed. In this case, we have no guarantees about the program  $C^*$  output by  $\mathcal{A}$  and we analyze two cases. One case is that  $C^*$  does not  $\varepsilon/2$ -approximate  $\tilde{C}$ . In this case, the probability that Test passes (outputs 1) is negligible. The other case is that  $C^*$  does  $\varepsilon/2$ -approximate  $\tilde{C}$ . In this case, since  $\mathbf{x}^* = (x_1^*, \dots, x_\ell^*)$  are  $\ell$  random values independent of  $\tilde{C}, C^*$ , it holds that

$$\Pr[\forall i \tilde{C}(x_i^*) \neq C^*(x_i^*)] = \prod_i \Pr[\tilde{C}(x_i^*) \neq C^*(x_i^*)] \leq (1 - \varepsilon/2)^\ell \leq \text{negl}(\lambda).$$

Therefore, if the distribution is a random the probability that  $\mathcal{D}$  outputs 1 is at most  $\text{negl}(\lambda)$ .

This means that  $\mathcal{D}$  has distinguishing advantage  $\delta - \text{negl}(\lambda)$ . □

*Proof of Lemma 4.2.* We define a sequence of hybrid games to prove the lemma.

**Hyb<sup>0</sup>:** This is the real game. The distribution is as follows. For  $K \leftarrow \text{Key}(1^\lambda)$ ,  $\mathbf{x} \leftarrow \mathcal{D}^\ell$  and  $\mathbf{y} \leftarrow \mathcal{R}^\ell$ ,

$$(\mathbf{x}, i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])).$$

**Hyb<sup>1</sup>:** In this game, we replace the pPRF key with a punctured one. The distribution is as follows. For  $K \leftarrow \text{Key}(1^\lambda)$ ,  $\mathbf{x} \leftarrow \mathcal{D}^\ell$ ,  $\mathbf{y} \leftarrow \mathcal{R}^\ell$ ,  $K\{\mathbf{x}\} := \text{Puncture}(K, \mathbf{x})$ , program  $C[K\{\mathbf{x}\}, \mathbf{x}, \mathbf{y}]$  in Figure 3,

$$(\mathbf{x}, i\mathcal{O}(C[K\{\mathbf{x}\}, \mathbf{x}, \mathbf{y}])).$$

**Constants:**  $x \in D^\ell$ ,  $y \in R^\ell$ , pPRF key  $K\{x\}$ .

**Inputs:**  $x \in D$

**Code:** If there exists  $i$  such that  $x = x_i$ , then output  $y_i$ . Else, output  $F(K\{x\}, x)$ .

Figure 3: Program  $C[K\{x\}, x, y]$

It holds that  $\text{Hyb}^0 \stackrel{c}{\approx} \text{Hyb}^1$  if  $i\mathcal{O}$  is an indistinguishability obfuscation. In  $\text{Hyb}^1$ , the PRF key is punctured, but program  $C[K\{x\}, x, y]$  does not compute  $F$  at points  $x$  due to the if-branch code. Thus, for all input  $x$ , the behavior of two programs is exactly the same and this is indistinguishable by  $i\mathcal{O}$ .

**Hyb<sup>2</sup>:** In this game, we replace  $y$  with outputs of pPRF. The distribution is as follows. For  $K \leftarrow \text{Key}(1^\lambda)$ ,  $x \leftarrow D^\ell$ ,  $y_i := F(K, x_i)$  for all  $i \in [\ell]$ ,  $K\{x\} := \text{Puncture}(K, x)$ , and program remains  $C[K\{x\}, x, y]$  shown in Figure 3.

It holds that  $\text{Hyb}^1 \stackrel{c}{\approx} \text{Hyb}^2$  by the security of the pPRF  $\mathcal{F}$ . The only difference is the distribution of  $y$ , but this is indistinguishable due to the pseudorandomness at punctured points of pPRF.

**Hyb<sup>3</sup>:** We replace the punctured PRF key with a standard one and erase the “if-branch” in the code of the obfuscated program. The distribution is as follows. For  $x \leftarrow D^\ell$ ,  $K \leftarrow \text{Key}(1^\lambda)$ , and program  $C[K]$  in Figure 4,

$$(x, i\mathcal{O}(C[K])).$$

It holds that  $\text{Hyb}^2 \stackrel{c}{\approx} \text{Hyb}^3$  if  $i\mathcal{O}$  is an indistinguishability obfuscation. In these two hybrid games, the

**Constants:** pPRF key  $K$ .

**Inputs:**  $x \in D$

**Code:** Output  $F(K, x)$ .

Figure 4: Program  $C[K]$

output of two programs is  $F(K, x)$  for any input  $x$ . Thus, the behavior of two programs are exactly the same due to the functionality preserving property under puncturing.

Note that in  $\text{Hyb}^3$  the vector  $x$  does not appear in the obfuscated program and therefore by “renaming” this vector we can write the distribution as  $x' \leftarrow D^\ell$ ,  $K \leftarrow \text{Key}(1^\lambda)$ , and program  $C[K]$  in Figure 4,

$$(x', i\mathcal{O}(C[K])).$$

**Hyb<sup>4</sup>:** We change the PRF key again, that is, we use  $K\{x\} := \text{Puncture}(K, x)$  and set  $y_i := F(K, x_i)$  for all  $i \in [\ell]$ . The distribution is as follows. For  $K \leftarrow \text{Key}(1^\lambda)$ ,  $x, x' \leftarrow D^\ell$ ,  $y_i := F(K, x_i)$  for all  $i \in [\ell]$ ,  $K\{x\} := \text{Puncture}(K, x)$ , and program  $C[K\{x\}, x, y]$  in Figure 3,

$$(x', i\mathcal{O}(C[K\{x\}, x, y])).$$

It holds that  $\text{Hyb}^3 \stackrel{c}{\approx} \text{Hyb}^4$  by the security of  $i\mathcal{O}$  (same as  $\text{Hyb}^2 \stackrel{c}{\approx} \text{Hyb}^3$ ).

Hyb<sup>5</sup>: We choose uniformly random  $\mathbf{y} \leftarrow \mathbb{R}^\ell$  again. For  $K \leftarrow \text{Key}^{1^\lambda}$ ,  $\mathbf{x}, \mathbf{x}' \leftarrow \mathbb{D}^\ell$ ,  $\mathbf{y} \leftarrow \mathbb{R}^\ell$ ,  $K\{\mathbf{x}\} := \text{Puncture}(K, \mathbf{x})$ , and the program in Figure 3,

$$(\mathbf{x}', i\mathcal{O}(C[K\{\mathbf{x}\}, \mathbf{x}, \mathbf{y}])).$$

It holds  $\text{Hyb}^4 \stackrel{c}{\approx} \text{Hyb}^5$  by the security of the pPRF (same as  $\text{Hyb}^1 \stackrel{c}{\approx} \text{Hyb}^2$ ).

Hyb<sup>6</sup>: We replace pPRF key  $K\{\mathbf{x}\}$  with a normal one. For  $K \leftarrow \text{Key}(1^\lambda)$ ,  $\mathbf{x}, \mathbf{x}' \leftarrow \mathbb{D}^\ell$ ,  $\mathbf{y} \leftarrow \mathbb{R}^\ell$ , and the program in Figure 1,

$$(\mathbf{x}', i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])).$$

It holds  $\text{Hyb}^5 \stackrel{c}{\approx} \text{Hyb}^6$  by the security of iO (same as  $\text{Hyb}^0 \stackrel{c}{\approx} \text{Hyb}^1$ ).

Thus, we obtained  $(\mathbf{x}, i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])) \stackrel{c}{\approx} (\mathbf{x}', i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}]))$ . □

## 4.2 Full Construction: Public-Key Detection and Security with Marking Oracle

We now construct a watermarking scheme with public-key detection and with security in the presence of a marking oracle (see Definition 3.4). We begin with a mark-embedding scheme (no message) that works for any approximation factor  $\varepsilon = 1/\text{poly}(\lambda)$  and then show how to convert it into a message-embedding scheme that works for approximation factors  $\varepsilon = \frac{1}{2} + 1/\text{poly}(\lambda)$ .

**Overview and Intuition.** There are several challenges that we need to overcome to go from the basic construction to our full construction. Not only does the basic construction crucially rely on the watermarking key  $wk$  being secret, it already fails to be secure even if the adversary has access to *either one of* the marking or detection oracles.

**Attack using the marking oracle:** Assume the adversary makes just one call to the marking oracle with an arbitrary key  $K'$  and gets a marked program  $\tilde{C}'$ . Then, on the marked-points  $x_1, \dots, x_\ell$  it will happen that  $\tilde{C}'(x_i) \neq F(K', x_i)$ , while for all other points equality holds. Therefore, the adversary can easily distinguish the marked-points from random points. In particular this leads to a concrete attack. Once the adversary gets a challenge marked program  $\tilde{C}$  for some random and unknown PRF key  $K$ , the adversary can create a program  $C^*$  which has  $K', \tilde{C}', \tilde{C}$  hard-coded and, on input  $x$ , if  $\tilde{C}'(x) \neq F(K', x)$  then  $C^*(x)$  outputs  $\perp$  else it outputs  $\tilde{C}(x)$ . The program  $C^*$  gives the same output as  $\tilde{C}$  on everything *other than* the marked-points and therefore very closely approximates  $\tilde{C}$  but will not be detected as marked with our detection procedure.

This attack shows that the set of marked-points cannot be static but must depend on the key being marked. But this raises a significant challenge since the detection procedure must be able to find marked-points and test the given program on them without knowing what the underlying key is.

**Attack using the detection oracle:** If the adversary gets the challenge marked program  $\tilde{C} = i\mathcal{O}(C[K, \mathbf{x}, \mathbf{y}])$ , he can create programs  $C'[v](x)$  which evaluate  $\tilde{C}(x)$  unless the first  $|v|$  bits of  $x$  match the value  $v$  in which case  $C'[v](x)$  outputs  $\perp$ . By calling the detection oracle with such program  $C'$  (for appropriate choices of  $v$ ) and seeing if detection succeeds or fails the adversary can use a variant of binary search to learn all the values  $\mathbf{x}$ . Once the adversary learns  $\mathbf{x}$  he can create a program  $C^*$  which evaluates  $\tilde{C}$  on all inputs other than the ones in  $\mathbf{x}$  and this will remove the mark while maintaining approximate correctness.

We now outline the main ideas for how to solve the above problems.

**Security with a Public-Detection Key:** Firstly, our main idea to get security in the presence of a detection oracle, is to use *super polynomially many* marked-points. More concretely, the set of marked-points is some super-polynomial sized pseudorandom subset  $S \subseteq D$  of negligible density. The detection procedure only tests the given program on some small (polynomial) number of random points in  $S$ . This way, even if the adversary learns the points that were used by polynomially executions of the detection procedure, it does not learn too much about the overall set of marked-points  $S$ . Since the set of marked-points is now super-polynomial we cannot simply hard-code the outputs of the marked program at the marked-points. We fix this next while also achieving security against a marking oracle.

Secondly, our main idea to get a scheme with a public-detection key, is to obfuscate the main components of the detection procedure so as to hide the set  $S$  of marked-points.

**Protection against the marking oracle:** Achieving security in the presence of a marking oracle is the most challenging part. To do this, we need to make sure that the set of marked-points  $x$  is different for each key  $K$  being marked while also ensuring that the detection procedure can find the marked-points without knowing  $K$ . We do this by introducing an additional concept that we call “find-points”  $\alpha$  which are static and do not depend on  $K$ . The set of marked-points  $x$  for some key  $K$  depends on the values  $F(K, \alpha)$  for the find-points  $\alpha$ . The detection procedure first tests the program on a find-point  $\alpha$  to find the corresponding marked-point  $x$  and only then tests if the value of the program on  $x$  is some special marked output. By calling the marking oracle on various keys  $K^{(1)}, K^{(2)}, \dots$  the adversary can learn something about the marked-points for each of these keys, but will not learn anything about the set of find-points or the set of marked-points for the challenge key  $K$ .

**Scheme Outline.** To make the above precise, assume we want to mark a PRF family with domain  $D = \{0, 1\}^n$  and range  $R = \{0, 1\}^m$ . We write elements in the domain  $D$  as  $x = x^L \| x^R$  where  $x^L \in \{0, 1\}^{n'}$  and  $x^R \in \{0, 1\}^{n-n'}$  for some prefix length  $n'$ . Our construction relies on three auxiliary PRFs  $F(K_1, \cdot)$ ,  $F(K_2, \cdot)$ ,  $F(K_3, \cdot)$ , with different domains/ranges as will be clear from context, which together make up the marking key  $mk = (K_1, K_2, K_3)$ .

For each value  $x^L \in \{0, 1\}^{n'}$  there is a corresponding *find-point*  $\alpha = F(K_1, x^L)$ , which is statically determined by the marking key. There is also a corresponding *marked-point*  $x = x^L \| x^R$  which depends on the output of the PRF  $F(K, \alpha)$  for the key  $K$  being marked and is computed by setting:  $\beta := F(K, \alpha)$ ,  $x^R := F(K_2, x^L \| \beta)$ . On a marked-point  $x$  as above, the marked program outputs a *marked-output*  $y = F(K_3, x^L \| \beta)$ . Note that there are  $\tau = 2^{n'}$  marked-points in total and, by picking  $n'$  appropriately, we can ensure that this is a super-polynomial number of marked-points but a negligible fraction of the total domain.

To mark a key  $K$  we obfuscate a program that gets some input  $x$ , tests if  $x$  is a marked point for  $K$ , and if so outputs a marked output  $y$ , and otherwise computes  $F(K, x)$ . In the symmetric-key setting, we would use a detection procedure that tests a program  $C'$  as follows. It first picks a random value  $x^L \leftarrow \{0, 1\}^{n'}$  and computes the corresponding find-point  $\alpha := F(K_1, x^L)$ . Then it computes  $\beta = C'(\alpha)$  and uses this to find the corresponding marked-point  $x = x^L \| x^R$  where  $x^R := F(K_2, x^L \| \beta)$ . Finally, it computes  $y = C'(x)$  and tests if  $y = F(K_3, x^L \| \beta)$  is a marked output (we can repeat this test several time to amplify). In the public-key detection setting, we essentially obfuscate the main components of the detection procedure and release the obfuscation as a public-detection key.

On an intuitive level, the above ensures that even if we are given the marked-points corresponding to some keys  $K'$ , we do not learn anything about the find-points or about the marked-points for the challenge key  $K$ . However, to remove a mark from the challenge marked program, an adversary would need to change the behavior of this program on a large fraction of either the find points or the marked points.

### 4.2.1 A Mark-Embedding Construction

We now give a detailed description of our mark-embedding scheme (no message) in the public-detection key setting and with security in the presence of a marking oracle.

**Setup, Parameters.** Recall that our goal is to construct a watermarking scheme for a pPRF family  $\mathcal{F}$  with  $F(K, \cdot) : \{0, 1\}^n \rightarrow \{0, 1\}^m$ . We assume  $\mathcal{F}$  has the injective property explained in Section 2.2. As in Section 4.1, we use notation  $D = \{0, 1\}^n$  and  $R = \{0, 1\}^m$ . (Note that  $n, m$  are polynomials in  $\lambda$  but we write  $n$  instead of  $n(\lambda)$  to simplify notation.)

Let  $n'' \leq n' \leq n$  be polynomials such that  $n - n' = \omega(\log \lambda)$ ,  $n' - n'' = \omega(\log \lambda)$  and  $n'' = \omega(\log \lambda)$ . Our construction will rely on three auxiliary pPRF families  $\mathcal{F}_1$  with  $F(K_1, \cdot) : \{0, 1\}^{n'} \rightarrow \{0, 1\}^n$ ,  $\mathcal{F}_2$  with  $F(K_2, \cdot) : \{0, 1\}^{n'+m} \rightarrow \{0, 1\}^{n-n'}$ , and  $\mathcal{F}_3$  with  $F(K_3, \cdot) : \{0, 1\}^{n'+m} \rightarrow \{0, 1\}^n$  where  $\mathcal{F}_2$  and  $\mathcal{F}_3$  are prefix pPRF families as in Definition 2.5. Our construction will also rely on a pseudorandom number generator (PRG)  $G(\cdot) : \{0, 1\}^{n''} \rightarrow \{0, 1\}^{n'}$ . In addition to the pPRF families  $\mathcal{F}, \mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3$  and the PRG  $G$ , we also rely on an indistinguishability obfuscator  $i\mathcal{O}$ .

Our proof will consist of  $O(\tau)$  hybrids where  $\tau := 2^{n'}$  is super-polynomial. In each of these hybrids we will rely on the security of  $\mathcal{F}, \mathcal{F}_3$  or  $i\mathcal{O}$ . This requires us to rely on complexity leveraging, where we assume that a PPT attacker cannot break the security of the schemes  $\mathcal{F}, \mathcal{F}_3$  or  $i\mathcal{O}$  with probability better than  $\mu = \mu(\lambda)$  such that  $\tau \cdot \mu$  remains negligible. This can be done if we (e.g.,) assume that  $\mathcal{F}, \mathcal{F}_3$  and the  $i\mathcal{O}$  have sub-exponential security. We do not need sub-exponential security for  $\mathcal{F}_1, \mathcal{F}_2$ , and  $G$ .

For a string  $x \in D$ , we will parse it as  $x = x^L \| x^R$  where  $x^L \in \{0, 1\}^{n'}$  and  $x^R \in \{0, 1\}^{n-n'}$ .

**Construction.** For any inverse-polynomial approximation factor  $\varepsilon = \varepsilon(\lambda)$  we set  $\ell = \ell(\lambda) = \lambda/\varepsilon^2$  and define our construction as follows.

**Gen( $1^\lambda$ ):** Choose the marking key  $mk := (K_1, K_2, K_3)$  consisting of randomly chosen keys for the PRFs  $\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3$  respectively. Output  $(dk, mk)$  where the detection key  $dk := (\tilde{D}_1, \tilde{D}_2)$  consists of two obfuscated programs  $\tilde{D}_1 \leftarrow i\mathcal{O}(D_1[K_1])$  and  $\tilde{D}_2 \leftarrow i\mathcal{O}(D_2[K_2, K_3])$  defined as:

$$D_1[K_1](r) := F(K_1, G(r)) \quad , \quad D_2[K_2, K_3](r, \beta) := (F(K_2, G(r) \| \beta), F(K_3, G(r) \| \beta)).$$

**Mark( $mk, K$ ):** Output  $\tilde{C} \leftarrow i\mathcal{O}(C[K, K_1, K_2, K_3])$  where the program  $C[K, K_1, K_2, K_3]$  is described in Figure 5.

**Constants:** pPRF keys  $K, K_1, K_2, K_3$ .

**Inputs:**  $x \in D$

**Code:** Parse  $x = x^L \| x^R$ . Set  $\alpha := F(K_1, x^L)$ ,  $\beta := F(K, \alpha)$ ,  $\hat{x}^R := F(K_2, x^L \| \beta)$ .  
If  $x^R = \hat{x}^R$  then output  $F(K_3, x^L \| \beta)$  else output  $F(K, x)$ .

Figure 5: Program  $C[K, K_1, K_2, K_3]$

**Detect( $dk, C'$ ):** For  $i \in [\ell]$ : select  $r_i \leftarrow \{0, 1\}^{n''}$  at random, set  $x_i^L := G(r_i)$ ,  $\alpha_i := \tilde{D}_1(r_i)$ ,  $\beta_i := C'(\alpha_i)$  and  $(x_i^R, y_i) := \tilde{D}_2(r_i, \beta_i)$ . If there is some  $i \in [\ell]$  such that  $C'(x_i^L \| x_i^R) = y_i$  then output marked, else output unmarked.

We can easily verify that the construction above satisfies statistical correctness, detection correctness, and meaningful properties.



**Security Proofs for the Mark-Embedding Construction.** We prove the non-removability security of the mark-embedding construction.

**Theorem 4.4** *The above watermarking construction is a secure mark-embedding scheme with a public-key detection and security in the presence of a marking oracle, assuming that:  $\mathcal{F}, \mathcal{F}_1$  are pPRFs, and  $\mathcal{F}_2, \mathcal{F}_3$  are prefix pPRFs,  $\mathcal{G}$  is a PRG and  $i\mathcal{O}$  is a indistinguishability obfuscator and  $\mathcal{F}, \mathcal{F}_3, i\mathcal{O}$  are sub-exponentially secure with distinguishing advantage  $\mu$  such that  $\tau\mu$  is negligible where  $\tau = 2^{n'}$ .*

**Lemma 4.5** *Under the same conditions as in Theorem 4.4, for all PPT distinguishers  $\mathcal{D}$  it holds that*

$$\left| \Pr[\mathcal{D}^{\mathcal{MO}(mk, \cdot)}(\tilde{C}, dk, \alpha, \mathbf{x}) = 1] - \Pr[\mathcal{D}^{\mathcal{MO}(mk, \cdot)}(\tilde{C}, dk, \alpha', \mathbf{x}') = 1] \right| \leq \text{negl}(\lambda).$$

where  $(mk = (K_1, K_2, K_3), dk) \leftarrow \text{Gen}(1^\lambda)$ ,  $K$  is a random key for  $\mathcal{F}$ ,  $\tilde{C} \leftarrow \text{Mark}(mk, K)$  and  $\mathcal{MO}(mk, \cdot)$  is the marking oracle. The values  $\alpha = (\alpha_1, \dots, \alpha_\ell)$  and  $\mathbf{x} = (x_1, \dots, x_\ell)$  correspond to random find-points and mark-points as sampled by the detection procedure and are chosen as follows. For  $i = 1, \dots, \ell$ : choose  $r_i \leftarrow \{0, 1\}^{n''}$  at random,  $x_i^L := \mathcal{G}(r_i)$ ,  $\alpha_i := F(K_1, x_i^L)$ ,  $\beta_i := F(K, \alpha_i)$ ,  $x_i^R := F(K_2, x_i^L \parallel \beta_i)$  and  $x_i := x_i^L \parallel x_i^R$ . On the other hand  $\alpha' \leftarrow \mathcal{D}^\ell$ ,  $\mathbf{x}' \leftarrow \mathcal{D}^\ell$  are uniformly random.

**Lemma 4.6** *If a PPT adversary  $\mathcal{A}$  breaks non-removability security of our construction with probability  $\delta(\lambda)$ , then we can construct a PPT distinguisher  $\mathcal{D}^{\mathcal{MO}(mk, \cdot)}$  such that*

$$\left| \Pr[\mathcal{D}^{\mathcal{MO}(mk, \cdot)}(\tilde{C}, dk, \alpha, \mathbf{x}) = 1] - \Pr[\mathcal{D}^{\mathcal{MO}(mk, \cdot)}(\tilde{C}, dk, \alpha', \mathbf{x}') = 1] \right| \geq \delta(\lambda) - \text{negl}(\lambda).$$

where the distributions are defined as in 4.5.

Lemma 4.5 and 4.6 immediately imply Theorem 4.4.

*Proof of Lemma 4.6.* Given a distribution  $(\tilde{C}, dk, \alpha^*, \mathbf{x}^*)$ , we construct a distinguisher  $\mathcal{D}$  as follows.  $\mathcal{D}$  gives  $dk$  and  $\tilde{C}$  to the adversary  $\mathcal{A}$  as a detection key and a challenge marked program, respectively. If  $\mathcal{A}$  attempts to query a PRF key  $K^{(j)}$  to the marking oracle, then  $\mathcal{D}$  passes it to  $\mathcal{MO}(mk, \cdot)$  and replies with the answer of  $\mathcal{MO}(mk, K^{(j)})$ . Eventually,  $\mathcal{A}$  outputs a program  $C^*$  and  $\mathcal{D}$  proceeds as follows:

1.  $\mathcal{D}$  runs  $\text{Test}(C^*, \tilde{C})$ . If  $\text{Test}$  outputs 0, then  $\mathcal{D}$  outputs 0.  $\text{Test}$  is described in Figure 2.
2. If for all  $i \in [\ell]$ ,  $(\tilde{C}(x_i^*) \neq C^*(x_i^*) \vee \tilde{C}(\alpha_i^*) \neq C^*(\alpha_i^*))$ ,  $\mathcal{D}$  outputs 1, otherwise 0.

The analysis of the distinguishing advantage is similar to that in the proof of Lemma 4.3. In particular, the analysis of  $\text{Test}$  is the same and ensures that if  $C^*$  does  $\varepsilon$ -approximate  $\tilde{C}$  then the test outputs 1 with overwhelming probability and if  $C^*$  does not  $\varepsilon/2$ -approximate  $\tilde{C}$  then the test outputs 0 with overwhelming probability. We now continue to the analysis of  $\mathcal{D}$ .

**Analysis of  $\mathcal{D}$ .** If the given distribution is the real one  $(\tilde{C}, dk, \alpha, \mathbf{x})$ , then  $\mathcal{A}$  wins the non-removability game with probability  $\delta$ , meaning that it outputs a program  $C^*$  that  $\varepsilon$ -approximates  $\tilde{C}$  but is not detected as marked by an execution of the detection procedure which evaluates  $C^*$  on the points  $\mathbf{x}, \alpha$ . If this happens, then  $\text{Test}$  passes and outputs 1 with overwhelming probability  $1 - \text{negl}(\lambda)$  and, for all  $i$ ,  $(C^*(x_i) \neq \tilde{C}(x_i) \vee \tilde{C}(\alpha_i^*) \neq C^*(\alpha_i^*))$ . Therefore, in this case,  $\mathcal{D}$  outputs 1 with probability at least  $\delta - \text{negl}(\lambda)$ .

If the given distribution is random one  $(\tilde{C}, dk, \alpha', \mathbf{x}')$ , then  $\mathbf{x}'$  and  $\alpha'$  are uniformly random and independently distributed. We analyze two cases. One case is that  $C^*$  does not  $\varepsilon/2$ -approximate  $\tilde{C}$ . In this case,

the probability that Test passes (outputs 1) is negligible. The other case is that  $C^*$  does  $\varepsilon/2$ -approximate  $\tilde{C}$ . In this case, since  $x'_1, \dots, x'_\ell, \alpha'_1, \dots, \alpha'_\ell$  are  $2\ell$  random values independent of  $\tilde{C}, C^*$ , it holds that

$$\begin{aligned} & \Pr \left[ \forall i \left( \tilde{C}(x'_i) \neq C^*(x'_i) \right) \vee \left( \tilde{C}(\alpha'_i) \neq C^*(\alpha'_i) \right) \right] \\ &= \prod_i \left( 1 - \Pr \left[ \left( \tilde{C}(x'_i) = C^*(x'_i) \right) \wedge \left( \tilde{C}(\alpha'_i) = C^*(\alpha'_i) \right) \right] \right) \\ &\leq (1 - (\varepsilon/2)^2)^\ell \leq \text{negl}(\lambda) \end{aligned}$$

if  $\ell = O(\lambda/\varepsilon^2)$ . Therefore, if the distribution is a random the probability that  $\mathcal{D}$  outputs 1 is at most  $\text{negl}(\lambda)$ . This means that  $\mathcal{D}$  has distinguishing advantage  $\delta - \text{negl}(\lambda)$ .  $\square$

*Proof of Lemma 4.5.* We define a sequence of hybrid games to prove this lemma. We call the left-hand side distribution in Lemma 4.5 the real distribution REAL and the right-hand side distribution is the random distribution RAND. Similarly in the proof of Lemma 4.2 in Section 4.1, we will transform the challenge program in the real distribution into another program that does not depend on  $K_1, K_2, K_3$  and then replace  $\mathbf{x}$  and  $\alpha$  with uniformly and independently random  $\mathbf{x}'$  and  $\alpha'$ , respectively. We define the following main hybrid games (but proving the indistinguishability of some of these will require more intermediate hybrids). See Table 1 for an abbreviated overview of the hybrids emphasizing the differences between them.

REAL: This is the real distribution (left-hand side in Lemma 4.5). The distinguisher is given the values

$$(\tilde{C}, dk = (\tilde{D}_1, \tilde{D}_2), \alpha, \mathbf{x})$$

and access to the marking oracle  $\mathcal{MO}(mk, \cdot)$ . Recall that  $(mk = (K_1, K_2, K_3), dk) \leftarrow \text{Gen}(1^\lambda)$ ,  $K$  is a random key for  $\mathcal{F}$  and  $\tilde{C} \leftarrow \text{Mark}(mk, K)$ . The values  $\alpha = (\alpha_1, \dots, \alpha_\ell)$  and  $\mathbf{x} = (x_1, \dots, x_\ell)$  are chosen as follows. For  $i = 1, \dots, \ell$ : choose  $r_i \leftarrow \{0, 1\}^{n'}$  at random,  $x_i^L := G(r_i)$ ,  $\alpha_i := F(K_1, x_i^L)$ ,  $\beta_i := F(K, \alpha_i)$ ,  $x_i^R := F(K_2, x_i^L \parallel \beta_i)$  and  $x_i := x_i^L \parallel x_i^R$ .

Hyb<sub>1</sub>: For  $i \in [\ell]$ , instead of choosing  $x_i^L := G(r_i)$ , we now choose  $x_i^L \leftarrow \{0, 1\}^{n'}$  uniformly at random.

Hyb<sub>2</sub>: Instead of choosing challenge program  $\tilde{C} \leftarrow i\mathcal{O}(C[K, K_1, K_2, K_3])$  where the program  $C[K, K_1, K_2, K_3]$  is described in Figure 5, we now set  $\tilde{C} \leftarrow i\mathcal{O}(C[K])$  where  $C[K](x) := F(K, x)$  just evaluates the PRF. Also, instead of choosing the detection key component  $\tilde{D}_2 \leftarrow i\mathcal{O}(D_2[K_2, K_3])$  as defined in the scheme description, we now choose  $\tilde{D}_2 \leftarrow i\mathcal{O}(D'_2[K, K_1, K_2, K_3])$  as defined in Figure 6. <sup>4</sup>

**Constants:** pPRF keys  $K, K_1, K_2, K_3$ .

**Inputs:**  $r, \beta$

**Code:** Set  $x^L := G(r)$ ,  $\alpha := F(K_1, x^L)$ ,  $x^R := F(K_2, x^L \parallel \beta)$ .

If  $\beta = F(K, \alpha)$ , then set  $y := F(K, x^L \parallel x^R)$  else set  $y := F(K_3, x^L \parallel \beta)$ . Output  $(x^R, y)$ .

Figure 6: Program  $D'_2[K, K_1, K_2, K_3]$

Hyb<sub>3</sub>: In this hybrid the values  $\mathbf{x}$  and  $\alpha$  are replaced with uniformly random elements  $\mathbf{x}' \leftarrow \mathcal{D}^\ell$  and  $\alpha' \leftarrow \mathcal{D}^\ell$ , respectively.

RAND: This is the random distribution (right-hand side in Lemma 4.5). The only changes from Hyb<sub>3</sub> are that the challenge program  $\tilde{C}$  and detection key  $\tilde{D}_2$  are changed back to the original programs but the values  $\mathbf{x}', \alpha'$  remain random.

Table 1: An overview of hybrid games

Hybrid game	Challenge: $i\mathcal{O}(\cdot)$	$dk: i\mathcal{O}(\cdot)$	$\alpha$	$x^L$	$x^R$
REAL	$C[K, K_1, K_2, K_3]$	$D_2[K_2, K_3]$	$F(K_1, x^L)$	$G(r)$	$F(K_2, x^L    \beta)$
Hyb <sub>1</sub>	$C[K, K_1, K_2, K_3]$	$D_2[K_2, K_3]$	$F(K_1, x^L)$	random	$F(K_2, x^L    \beta)$
Hyb <sub>2</sub>	$C[K]$	$D'_2[K, K_1, K_2, K_3]$	$F(K_1, x^L)$	random	$F(K_2, x^L    \beta)$
Hyb <sub>3</sub>	$C[K]$	$D'_2[K, K_1, K_2, K_3]$	random	random	random
RAND	$C[K, K_1, K_2, K_3]$	$D_2[K_2, K_3]$	random	random	random

**Lemma 4.7** *If  $G$  is a pseudo-random number generator, then  $\text{REAL} \stackrel{c}{\approx} \text{Hyb}_1$ .*

*Proof of Lemma 4.7.* These hybrid games are indistinguishable due to the pseudorandomness of  $G$ .  $\square$

**Notation.** Before we proceed to the rest of the proof, we introduce some notation for puncturing PRFs. We write  $K\{x \rightarrow y\}$  to denote a PRF key  $K$  which is punctured at the points  $x$  with a *hard-coded output*  $y$ . In other words this denotes the process of first creating a punctured key  $K\{x\} := \text{Puncture}(K, x)$  and storing the tuple  $K\{x \rightarrow y\} = (K\{x\}, x, y)$ . We define  $F(K\{x \rightarrow y\}, x) = y$  and else  $F(K\{x \rightarrow y\}, x') = F(K\{x\}, x') = F(K, x')$  when  $x' \neq x$ .

We also use similar notation for prefix pPRFs where we write  $K\{s||* \rightarrow y\} = (K\{s||*\}, s, y)$ . For any input  $x$  that has  $s$  as a prefix, we define  $F(K\{s||* \rightarrow y\}, x) = y$  and for all other  $x'$  we define  $F(K\{s||* \rightarrow y\}, x') = F(K\{s||*\}, x') = F(K, x')$ .

Lastly we use the notation  $K\{x \rightarrow \perp\}$  to denote  $K\{x \rightarrow y_\perp\}$  for some fixed “dummy value”  $y_\perp$  in the range of  $F(K, \cdot)$ , say the all 0 string.

**Lemma 4.8** *If  $\mathcal{F}$  and  $\mathcal{F}_3$  are sub-exponentially secure pPRFs,  $\mathcal{F}$  satisfies the injective property in Definition 2.6, and  $i\mathcal{O}$  is a sub-exponentially secure indistinguishability obfuscation, so that the distinguishing advantage for each of these is bounded by some function  $\mu$  such that  $\tau\mu = \text{negl}(\lambda)$ , then  $\text{Hyb}_1 \stackrel{c}{\approx} \text{Hyb}_2$ .*

*Proof of Lemma 4.8.* To prove indistinguishability of  $\text{Hyb}_1 \stackrel{c}{\approx} \text{Hyb}_2$ , we introduce intermediate hybrid games corresponding to each one of the values  $x^L \in \{0, 1\}^{n'}$ . We interchangeably identify the strings  $x^L \in \{0, 1\}^{n'}$  with integers  $j \in \{0, \dots, \tau - 1\}$  for  $\tau = 2^{n'}$  via the usual binary representation. We define the following hybrids games (see Table 2 for a concise overview).

$\text{Hyb}_1^j$ : Everything is the same as in  $\text{Hyb}_1$  except that the challenge program  $\tilde{C}$  and the detection key component  $\tilde{D}_2$  are changed to:

- $\tilde{C} \leftarrow i\mathcal{O}(C^j[K, K_1, K_2, K_3])$  where  $C^j[K, K_1, K_2, K_3]$  is described in Figure 7,
- $\tilde{D}_2 \leftarrow i\mathcal{O}(D_2^j[K, K_2, K_3])$  where  $D_2^j[K, K_2, K_3]$  is described in Figure 8.

$\text{Hyb}_1^{j-1}$ : Everything is the same as in  $\text{Hyb}_1^j$  except that we always use the punctured keys

$$K\{(j||x_j^R) \rightarrow \perp\} \quad \text{and} \quad K_3\{(j||\beta_j) \rightarrow y_j\}$$

in all obfuscated programs in place of  $K, K_3$ , where we define  $\alpha_j = F(K_1, j)$ ,  $\beta_j = F(K, \alpha_j)$ ,  $x_j^R = F(K_2, j||\beta_j)$  and  $y_j = F(K_3, j||\beta_j)$ . In particular, this means that we obfuscate the following programs:

<sup>4</sup>As a sanity check, note that this ensures that the detection procedure still outputs marked when executed on  $\tilde{C}$ .

**Constants:** pPRF keys  $K, K_1, K_2, K_3$ .

**Inputs:**  $x \in \mathcal{D}$

**Code:** Parse  $x = x^L \| x^R$ . Set  $\alpha := F(K_1, x^L)$ ,  $\beta := F(K, \alpha)$ ,  $\hat{x}^R := F(K_2, x^L \| \beta)$ .  
If  $x^R \neq \hat{x}^R$  or  $x^L < j$ , then output  $F(K, x)$  else output  $F(K_3, x^L \| \beta)$ .

Figure 7: Program  $C^j[K, K_1, K_2, K_3]$

**Constants:** pPRF keys  $K, K_1, K_2, K_3$ .

**Inputs:**  $r, \beta$

**Code:** Set  $x^L := G(r)$ ,  $\alpha := F(K_1, x^L)$ ,  $x^R := F(K_2, x^L \| \beta)$ .  
If  $x^L < j$  and  $\beta = F(K, \alpha)$ , then set  $y := F(K, x^L \| x^R)$  else set  $y := F(K_3, x^L \| \beta)$ .  
Output  $(x^R, y)$ .

Figure 8: Program  $D_2^j[K, K_1, K_2, K_3]$

Table 2: An overview of hybrid games from  $\text{Hyb}_1$  to  $\text{Hyb}_2$

Hybrid	Challenge: $i\mathcal{O}(C[\cdot])$	$dk$ : $i\mathcal{O}(D_2[\cdot])$	$i$ -th query: $i\mathcal{O}(C[\cdot])$	hard-wired $y_j$
$\text{Hyb}_1^0$	$C[K, K_1, K_2, K_3]$	$D_2[K_2, K_3]$	$C[K^{(i)}, K_1, K_2, K_3]$	none
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\text{Hyb}_1^j$	$C^j[K, K_1, K_2, K_3]$	$D_2^j[K, K_1, K_2, K_3]$	$C[K^{(i)}, K_1, K_2, K_3]$	none
$\text{Hyb}_1^{j-1}$	$C^j[\text{punctured}]^*$	$D_2^j[\text{punctured}]^*$	$C[\text{punctured}]^*$	$F(K_3, j \  \beta)$
$\text{Hyb}_1^{j-2}$	$C^j[\text{punctured}]^*$	$D_2^j[\text{punctured}]^*$	$C[\text{punctured}]^*$	$F(K, j \  x_j^R)$
$\text{Hyb}_1^{j+1}$	$C^{j+1}[K, K_1, K_2, K_3]$	$D_2^{j+1}[K, K_1, K_2, K_3]$	$C[K^{(i)}, K_1, K_2, K_3]$	none
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\text{Hyb}_2$	$C^{\tau+1}[K, K_1, K_2, K_3]$	$D_2^{\tau+1}[K, K_1, K_2, K_3]$	$C[K^{(i)}, K_1, K_2, K_3]$	none

\* In this table, [punctured] means keys  $K$  and  $K_3$  are punctured. Note that keys  $K_1, K_2$ , and  $K^{(i)}$  are *not* punctured. See the definitions of hybrid games for details.

- The challenge program  $\tilde{C}$  obfuscates  $C^j[K\{j\|x_j^R \rightarrow \perp\}, K_1, K_2, K_3\{j\|\beta_j \rightarrow y_j\}]$ .
- The detection key component  $\tilde{D}_2$  obfuscates  $D_2^j[K\{j\|x_j^R \rightarrow \perp\}, K_1, K_2, K_3\{j\|\beta_j \rightarrow y_j\}]$ .
- Each query to the marking oracle with some PRF key  $K^{(i)}$  outputs a program  $\tilde{C}^{(i)}$  that obfuscates  $C[K^{(i)}, K_1, K_2, K_3\{j\|\beta_j \rightarrow y_j\}]$ .

$\text{Hyb}_1^{j-2}$ : Everything is the same as in  $\text{Hyb}_1^{j-1}$  except that use the hard-wired value  $y_j := F(K, j \| x_j^R)$  instead of  $y_j = F(K_3, j \| \beta_j)$ .

Let  $\mu$  be the maximal distinguishing advantage of a PPT attacker against the iO scheme and the pPRFs  $\mathcal{F}, \mathcal{F}_3$ . Let  $q$  be the number of queries to the marking oracle. Then the following hybrids are indistinguishable for all  $j \in \{0, \dots, \tau\}$ .

**Claim:**  $\text{Hyb}_1^j \stackrel{c}{\approx}_{(q+2)\mu} \text{Hyb}_1^{j-1}$ . This holds by iO security since the programs being obfuscated in both hybrids are functionally equivalent. In particular, none of the programs ever evaluates  $F(K, j \| x_j^R)$  and

therefore we can use a punctured key  $K\{(j\|x_j^R) \rightarrow \perp\}$  in place of  $K$  without changing functionality. On the other hand the keys  $K_3$  and  $K_3\{(j\|\beta_j) \rightarrow y_j\}$  are functionally equivalent.

There are  $q+2$  obfuscated programs in total and therefore we pay a factor of  $q+2$  in the distinguishing advantage: the  $q$  outputs of the marking oracle  $\tilde{C}^{(i)}$ , the challenge program  $\tilde{C}$ , and the detection key program  $\tilde{D}_2$ .

**Claim:**  $\text{Hyb}_1^{j-1} \stackrel{c}{\approx}_{2\mu} \text{Hyb}_1^j$ -2. This holds by the pPRF security of  $\mathcal{F}$  and  $\mathcal{F}_3$ . In particular, given the punctured keys  $K\{j\|x_j^R\}$  and  $K_3\{j\|\beta_j\}$ , the value  $F(K_3, j\|\beta_j)$  is indistinguishable from a uniformly random value which is in turn indistinguishable from  $F(K, j\|x_j^R)$ .

**Claim:**  $\text{Hyb}_1^{j-2} \stackrel{c}{\approx}_{(q+2)\mu} \text{Hyb}_1^{j+1}$ . This holds by iO security since the programs being obfuscated in both hybrids are functionally equivalent. In particular, in both hybrids, the obfuscated challenge program  $\tilde{C}$  and detection-key program  $\tilde{D}_2$  are using  $F(K, j\|\beta_j)$  in place of  $F(K_3, j\|\beta_j)$ . On the other hand, each of the programs output by the marking oracle  $\tilde{C}^{(i)}$  never evaluates  $F(K_3, j\|\beta_j)$  and hence it does not matter whether these programs use a punctured key or not. This follows since, each program  $\tilde{C}^{(i)}$  evaluates  $F(K_3, j\|\beta_{i,j})$  where  $\beta_{i,j} = F(K^{(i)}, \alpha_j)$  is different from  $\beta_j = F(K, \alpha_j)$  by the injective property of the pPRF  $\mathcal{F}$ .<sup>5</sup>

Moreover, we have  $\text{Hyb}_1 \stackrel{c}{\approx}_{2\mu} \text{Hyb}_1^0$  and  $\text{Hyb}_1^{\ell} \stackrel{c}{\approx}_{2\mu} \text{Hyb}_2$  since the programs being obfuscated to get  $\tilde{C}$  and  $\tilde{D}_2$  in these hybrids are functionally equivalent. Combining this and the the above claims, we get  $\text{Hyb}_1 \stackrel{c}{\approx}_{O(q\tau\mu)} \text{Hyb}_2$ . Since  $\tau\mu = \text{negl}(\lambda)$  and  $q = \text{poly}(\lambda)$  we get the computational indistinguishability  $\text{Hyb}_1 \stackrel{c}{\approx} \text{Hyb}_2$  which proves the lemma.  $\square$

**Lemma 4.9** *If  $\mathcal{F}_1$  is a secure pPRF,  $\mathcal{F}_2, \mathcal{F}_3$  are secure prefix pPRFs, and iO is secure indistinguishability obfuscation, then it holds  $\text{Hyb}_2 \stackrel{c}{\approx} \text{Hyb}_3$ .*

*Proof of Lemma 4.9.* To prove the indistinguishability of these hybrid games, we introduce intermediate hybrid games for  $i = 1, \dots, \ell + 1$ .

**Hyb<sub>2</sub><sup>i</sup>:** This is the same as  $\text{Hyb}_2$  except for how the values  $\alpha = (\alpha_1, \dots, \alpha_\ell)$ ,  $\mathbf{x} = (x_1, \dots, x_\ell)$  are chosen. For  $i' < i$ , instead of choosing  $\alpha_{i'} = F(K_1, x_{i'}^L)$ ,  $\beta_{i'} = F(K, \alpha_{i'})$ ,  $x_{i'}^R = F(K_2, x_{i'}^L\|\beta_{i'})$  and setting  $x_{i'} = x_{i'}^L\|x_{i'}^R$  we now choose uniformly random values  $\alpha_{i'} \leftarrow \mathcal{D}$ ,  $x_{i'} \leftarrow \mathcal{D}$ .

**Hyb<sub>2</sub><sup>i</sup>-1:** This is the same as  $\text{Hyb}_2^i$  but instead of the PRF keys  $K_1$  and  $K_2, K_3$  we use punctured keys  $K_1\{x_i^L\}$ ,  $K_2\{x_i^L\|*\}$ ,  $K_3\{x_i^L\|*\}$  in all obfuscated programs – namely, in the detection key  $\tilde{D}_1, \tilde{D}_2$  and in the marked programs  $\tilde{C}^{(j)}$  output by the marking oracle (note that  $K_1, K_2, K_3$  do not appear in the challenge program  $\tilde{C}$  at this point). We do this as follows:

- For the detection key  $\tilde{D}_1, \tilde{D}_2$  we can simply output  $\perp$  on the punctured inputs. In other words,  $\tilde{D}_1, \tilde{D}_2$  are respectively obfuscations of

$$D_1[K_1\{x_i^L \rightarrow \perp\}] \quad , \quad D_2'[K, K_1\{x_i^L \rightarrow \perp\}, K_2\{x_i^L\|* \rightarrow \perp\}, K_3\{x_i^L\|* \rightarrow \perp\}].$$

<sup>5</sup>Technically, we need to ensure that the distinguisher never queries the marking oracle with the key  $K^{(i)} = K$  contained in the challenge program. This is without loss of generality since if the distinguisher can ever guess  $K$  then it can already perfectly distinguish the two challenge distributions without making additional calls to the marking oracle. In other words, we can convert any distinguisher into one with the above restriction while maintaining the same distinguishing advantage.

Table 3: An overview of hybrid games from  $\text{Hyb}_2^i$  to  $\text{Hyb}_2^{i+1}$ 

Hybrid	$K_1$	$K_b$ $b \in \{2, 3\}$	$\alpha_i$	$x_i^R$	hard-wired PRF outputs in $C[K^{(j)}, \cdot]$	hard-wired PRF outputs in $D_1[\cdot]$ and $D_2'[\cdot]$
$\text{Hyb}_2^i$	$K_1$	$K_b$	$F(K_1, x_i^L)$	$F(K_2, x_i^L \parallel \beta_i)$	none	none
$\text{Hyb}_2^{i-1}$	$K_1\{x_i^L\}$	$K_b\{x_i^L \parallel *\}$	$F(K_1, x_i^L)$	$F(K_2, x_i^L \parallel \beta_i)$	$\hat{x}_{j,i}^R$   $y_{j,i}$	$\perp$
$\text{Hyb}_2^{i-2}$	$K_1\{x_i^L\}$	$K_b\{x_i^L \parallel *\}$	$F(K_1, x_i^L)$	random	random   random	$\perp$
$\text{Hyb}_2^{i-3}$	$K_1\{x_i^L\}$	$K_b\{x_i^L \parallel *\}$	random	random	random   random	$\perp$
$\text{Hyb}_2^{i-4}$	$K_1\{x_i^L\}$	$K_b\{x_i^L \parallel *\}$	random	random	$\hat{x}_{j,i}^R$   $y_{j,i}$	$\perp$
$\text{Hyb}_2^{i+1}$	$K_1$	$K_b$	random	random	none	none

- For each query  $K^{(j)}$  to the marking oracle, we compute  $\alpha_i = F(K_1, x_i^L)$ ,  $\beta_{j,i} = F(K^{(j)}, \alpha_i)$ ,  $x_{j,i}^R = F(K_2, x_i^L \parallel \beta_{j,i})$ ,  $y_{j,i} = F(K_3, x_i^L \parallel \beta_{j,i})$  and we obfuscate the program

$$C[K^{(j)}, K_1\{x_i^L \rightarrow \perp\}, K_2\{x_i^L \parallel * \rightarrow x_{j,i}^R\}, K_3\{x_i^L \parallel * \rightarrow y_{j,i}\}].$$

We can rewrite the code of the above program in a simplified (but functionally equivalent) manner as shown Figure 9.

<p><b>Constants:</b> pPRF keys <math>K^{(j)}, K_1\{x_i^L \rightarrow \perp\}, K_2\{x_i^L \parallel * \rightarrow x_{j,i}^R\}, K_3\{x_i^L \parallel * \rightarrow y_{j,i}\}</math>.</p> <p><b>Inputs:</b> <math>x \in D</math></p> <p><b>Code:</b> Parse <math>x = x^L \parallel x^R</math>. If <math>x^L = x_i^L</math>, then set <math>\hat{x}^R := x_{j,i}^R</math> and <math>\hat{y} = y_{j,i}</math>.                  Else, set <math>\alpha := F(K_1\{x_i^L\}, x^L)</math>, <math>\beta := F(K^{(j)}, \alpha)</math>, <math>\hat{x}^R := F(K_2\{x_i^L \parallel *\}, x^L \parallel \beta)</math>, <math>\hat{y} = F(K_3, x^L \parallel \beta)</math>.                  If <math>x^R = \hat{x}^R</math>, then output <math>\hat{y}</math> else output <math>F(K^{(j)}, x)</math>.</p>
--

 Figure 9: Program  $C[K^{(j)}, K_1\{x_i^L \rightarrow \perp\}, K_2\{x_i^L \parallel * \rightarrow \hat{x}_{j,i}^R\}, K_3\{x_i^L \parallel * \rightarrow y_{j,i}\}]$ 

$\text{Hyb}_2^{i-2}$ : This is the same as  $\text{Hyb}_2^{i-1}$  except that:

- When answering the  $j$ 'th marking-oracle query, we choose the values  $x_{j,i}^R$  and  $y_{j,i}$  uniformly at random instead of computing  $\alpha_i = F(K_1, x_i^L)$ ,  $\beta_{j,i} = F(K^{(j)}, \alpha_i)$ ,  $x_{j,i}^R = F(K_2, x_i^L \parallel \beta_{j,i})$ ,  $y_{j,i} = F(K_3, x_i^L \parallel \beta_{j,i})$ .
- In  $x = (x_1, \dots, x_\ell)$  we now sample  $x_i = x_i^L \parallel x_i^R$  by choosing  $x_i^R$  uniformly at random instead of computing  $\beta_i = F(K, \alpha_i)$ ,  $x_i^R = F(K_2, x_i^L \parallel \beta_i)$ . This is equivalent to sampling  $x_i \leftarrow D$  uniformly at random.

$\text{Hyb}_2^{i-3}$ : This is the same as  $\text{Hyb}_2^{i-2}$  except that, in  $\alpha = (\alpha_1, \dots, \alpha_\ell)$ , we now select  $\alpha_i \leftarrow D$  uniformly at random instead of  $\alpha_i = F(K_1, x_i^L)$ .

$\text{Hyb}_2^{i-4}$ : This is the same as  $\text{Hyb}_2^{i-3}$  except that, when answering the  $j$ 'th marking-oracle query, we switch back to choosing  $x_{j,i}^R$  and  $y_{j,i}$  by computing  $\alpha_i = F(K_1, x_i^L)$ ,  $\beta_{j,i} = F(K^{(j)}, \alpha_i)$ ,  $x_{j,i}^R = F(K_2, x_i^L \parallel \beta_{j,i})$ ,  $y_{j,i} = F(K_3, x_i^L \parallel \beta_{j,i})$ . Note that the value  $\alpha_i$  in the vector  $\alpha$  is still chosen uniformly at random.

We describe an overview of hybrid games in Table 3 to emphasize differences. The following hybrids are indistinguishable for all  $i \in \{1, \dots, \ell + 1\}$ .

**Claim:**  $\text{Hyb}_2^i \stackrel{c}{\approx} \text{Hyb}_2^{i-1}$ . This follows by iO security since the programs being obfuscated are functionally equivalent. This is clear for the responses of the marking oracle. For the detection key components  $\tilde{D}_1, \tilde{D}_2$ , note that the programs  $D_1, D_2'$  only evaluate  $K_1, K_2, K_3$  on the points of the form  $x^\perp$  or  $x^\perp \| *$  where  $x^\perp = G(r)$  for some  $r$ . However, since  $x_i^\perp$  is chosen uniformly at random, it is not of this form (with overwhelming probability) and therefore we can puncture the keys at  $x_i^\perp$  or  $x_i^\perp \| *$  without changing functionality.

**Claim:**  $\text{Hyb}_2^{i-1} \stackrel{c}{\approx} \text{Hyb}_2^{i-2}$ . This follows by the prefix pPRF security of  $\mathcal{F}_2, \mathcal{F}_3$ . In particular, since we use punctured keys  $K_2\{x_i^\perp \| *\}, K_3\{x_i^\perp \| *\}$  their outputs are indistinguishable from random on inputs of the form  $x_i^\perp \| *$ . Note that by injectivity of  $\mathcal{F}$  the values  $\beta_i$  and  $\beta_{j,i}$  are distinct and therefore we can replace the PRF outputs by random and independent values.

**Claim:**  $\text{Hyb}_2^{i-2} \stackrel{c}{\approx} \text{Hyb}_2^{i-3}$ . This follows by the pPRF security of  $\mathcal{F}_1$ . In particular, since we use a punctured key  $K_1\{x_i^\perp\}$  we can replace  $F(K_1, x_i^\perp)$  by uniform.

**Claim:**  $\text{Hyb}_2^{i-3} \stackrel{c}{\approx} \text{Hyb}_2^{i-4}$ . This follows by the prefix pPRF security of  $\mathcal{F}_2, \mathcal{F}_3$  (same as showing  $\text{Hyb}_2^{i-1} \stackrel{c}{\approx} \text{Hyb}_2^{i-2}$ ).

**Claim:**  $\text{Hyb}_2^{i-4} \stackrel{c}{\approx} \text{Hyb}_2^{i+1}$ : This follows by iO security (same as showing  $\text{Hyb}_2^i \stackrel{c}{\approx} \text{Hyb}_2^{i-1}$ ).

Lastly, by noting that  $\text{Hyb}_2 = \text{Hyb}_2^1$  and  $\text{Hyb}_2^{\ell+1} = \text{Hyb}_3$  we get  $\text{Hyb}_2 \stackrel{c}{\approx} \text{Hyb}_3$  which proves the lemma.  $\square$

**Lemma 4.10**  $\text{Hyb}_3 \stackrel{c}{\approx} \text{RAND}$ .

*Proof of Lemma 4.10.* This proof is the same as the proof of Lemma 4.8 (in reverse manner).  $\square$

By Lemma 4.7, 4.8, 4.9, 4.10, we complete the proof of Lemma 4.5.  $\square$

## 4.2.2 A Message-Embedding Construction

In this section, we propose a message-embedding watermarking scheme in the public-key detection setting and with selective-message security in the presence of a marking oracle. We can get a scheme that allows for any approximation factor  $\varepsilon(\lambda) = \frac{1}{\sqrt{2}} + 1/\text{poly}(\lambda)$ .

**Setup.** Our construction can be thought of as an “upgrade” of the construction in Section 4.2.1. As before, our goal is to construct a watermarking scheme for a pPRF family  $\mathcal{F}$  with  $F(K, \cdot) : \{0, 1\}^n \rightarrow \{0, 1\}^m$ . We set the message space to be some  $\mathcal{M} \subseteq \{0, 1\}^m$ . Let the pPRF  $\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3$  and the PRG  $G$  be defined as in Section 4.2.1.

**Construction.** For any approximation factor  $\varepsilon(\lambda) = \frac{1}{\sqrt{2}} + \rho(\lambda)$  where  $\rho(\lambda)$  is some inverse polynomial, we set  $\ell = \ell(\lambda) = \lambda/\rho^2$  and define our construction as follows.

**Gen( $1^\lambda$ ):** This is exactly the same as the key-generation procedure of our mark-embedding construction. Choose the marking key  $mk := (K_1, K_2, K_3)$  consisting of randomly chosen keys for the PRFs  $\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3$  respectively. Output  $(dk, mk)$  where the detection key  $dk := (\tilde{D}_1, \tilde{D}_2)$  consists of two obfuscated programs  $\tilde{D}_1 \leftarrow i\mathcal{O}(D_1[K_1])$  and  $\tilde{D}_2 \leftarrow i\mathcal{O}(D_2[K_2, K_3])$  defined as:

$$D_1[K_1](r) := F(K_1, G(r)) \quad , \quad D_2[K_2, K_3](r, \beta) := (F(K_2, G(r) \| \beta), F(K_3, G(r) \| \beta)).$$

**Constants:** pPRF keys  $K, K_1, K_2, K_3$  and message  $\text{msg}$ .

**Inputs:**  $x \in \mathcal{D}$

**Code:** Parse  $x = x^L \| x^R$ . Set  $\alpha := F(K_1, x^L)$ ,  $\beta := F(K, \alpha)$ ,  $\hat{x}^R := F(K_2, x^L \| \beta)$ .  
If  $x^R = \hat{x}^R$  then output  $F(K_3, x^L \| \beta) \oplus \text{msg}$  else output  $F(K, x)$ .

Figure 10: Program  $C[K, K_1, K_2, K_3, \text{msg}]$

**Mark( $mk, K, \text{msg}$ ):** Output  $i\mathcal{O}(C[K, K_1, K_2, K_3, \text{msg}])$  where the program  $C[K, K_1, K_2, K_3, \text{msg}]$  is described in Figure 10.

**Detect( $dk, C'$ ):** For  $i = 1, \dots, \ell$ : Select  $r_i$  at random, set  $x_i^L := G(r_i)$ ,  $\alpha_i := \tilde{D}_1(r_i)$ ,  $\beta_i := C'(\alpha_i)$  and  $(x_i^R, y_i) := \tilde{D}_2(r_i \| \beta_i)$ , and compute  $\text{msg}_i := C'(x_i^L \| x_i^R) \oplus y_i$ . If there exists a “majority value”  $\text{msg}$  such that  $|\{i : \text{msg} = \text{msg}_i\}| > \ell/2$  then output  $\text{msg}$ , else output unmarked.

We can easily verify that the above construction above satisfies statistical correctness, detection correctness, and meaningful properties.

**Security Proofs for the Message-Embedding Construction.** We prove the non-removability security of the above construction in the selective-message setting, following essentially the same strategy as the mark-embedding construction with small modifications.

**Theorem 4.11** *The above construction is a message-embedding watermarking scheme with public-detection and selective-message security in the presence of a marking oracle for any approximation factor  $\varepsilon(\lambda) = \frac{1}{\sqrt{2}} + 1/\text{poly}(\lambda)$  under the same conditions as in Theorem 4.4.*

The above theorem follows from the following two lemmas.

**Lemma 4.12** *Under the conditions of Theorem 4.4, it holds that for all PPT  $\mathcal{D}$  and all message  $\text{msg} \in \mathcal{M}$ :*

$$\left| \Pr[\mathcal{D}^{\mathcal{MO}(mk, \cdot, \cdot)}(\tilde{C}, dk, \alpha, \mathbf{x}) = 1] - \Pr[\mathcal{D}^{\mathcal{MO}(mk, \cdot, \cdot)}(\tilde{C}, dk, \alpha', \mathbf{x}') = 1] \right| \leq \text{negl}(\lambda).$$

where  $(mk = (K_1, K_2, K_3), dk) \leftarrow \text{Gen}(1^\lambda)$ ,  $K$  is a random pPRF key,  $\tilde{C} \leftarrow \text{Mark}(mk, K, \text{msg})$ ,  $\mathcal{MO}(mk, \cdot, \cdot)$  is the marking oracle that gets as input a pPRF key  $K^{(j)}$  and a message  $\text{msg}_j$  and outputs  $\tilde{C} \leftarrow \text{Mark}(mk, K^{(j)}, \text{msg}_j)$ . The values  $\alpha = (\alpha_1, \dots, \alpha_\ell)$  and  $\mathbf{x} = (x_1, \dots, x_\ell)$  correspond to random find-points and mark-points as sampled by the detection procedure and are chosen as follows. For  $i = 1, \dots, \ell$ : choose  $r_i \leftarrow \{0, 1\}^{n''}$  at random,  $x_i^L := G(r_i)$ ,  $\alpha_i := F(K_1, x_i^L)$ ,  $\beta_i := F(K, \alpha_i)$ ,  $x_i^R := F(K_2, x_i^L \| \beta_i)$  and  $x_i := x_i^L \| x_i^R$ . On the other hand  $\alpha' \leftarrow \mathcal{D}^\ell$ ,  $\mathbf{x}' \leftarrow \mathcal{D}^\ell$  are uniformly random.

**Lemma 4.13** *If a PPT adversary  $A$  breaks selective-message non-removability security of the given scheme with probability  $\delta(\lambda)$  and approximation factor  $\varepsilon > \frac{1}{\sqrt{2}} + \rho$ , then we can construct a PPT distinguisher  $\mathcal{D}^{\mathcal{MO}(mk, \cdot, \cdot)}$  and some message  $\text{msg} \in \mathcal{M}$  such that*

$$\left| \Pr[\mathcal{D}^{\mathcal{MO}(mk, \cdot, \cdot)}(\tilde{C}, dk, \alpha, \mathbf{x}) = 1] - \Pr[\mathcal{D}^{\mathcal{MO}(mk, \cdot, \cdot)}(\tilde{C}, dk, \alpha', \mathbf{x}') = 1] \right| \geq \delta(\lambda) - \text{negl}(\lambda).$$

where the distributions are those of 4.12.



*Proof of Lemma 4.13.* This proof is similar to that of Lemmas 4.3 and 4.6 with some minor differences. In particular, the distinguisher  $\mathcal{D}$  sets  $\text{msg}$  to be the message chosen by the adversary  $\mathcal{A}$  at the beginning of the game, and gets the challenge values  $(\tilde{C}, dk, \alpha^*, x^*)$ , where  $\alpha^*$  and  $x^*$  is either  $(\alpha, x)$  or  $(\alpha', x')$ . It gives  $\tilde{C}, dk$  to  $\mathcal{A}$  and simulates the marking oracle for  $\mathcal{A}$  using its own oracle. At some point  $\mathcal{A}$  outputs  $C^*$ .

1.  $\mathcal{D}$  runs a test that finds an estimate  $\hat{\rho}$  for  $\Pr_{x \leftarrow \mathcal{D}}[\tilde{C}(x) = C^*(x)] - \frac{1}{\sqrt{2}}$  which is within  $\rho/4$  of the actual probability with overwhelming probability. (This can be done the same way as in Test described in Figure. 2 with appropriate changes in parameters). If  $\hat{\rho} < \frac{3}{4}\rho$  then output 0.
2. If for majority of  $i \in [\ell]$ ,  $(\tilde{C}(x_i^*) \neq C^*(x_i^*)) \vee (\tilde{C}(\alpha_i^*) \neq C^*(\alpha_i^*))$ ,  $\mathcal{D}$  outputs 1, otherwise 0.

If the distribution of  $(x^*, \alpha^*)$  is that of  $(x, \alpha)$  then with probability  $\delta$  the adversary  $\mathcal{A}$  outputs a “winning program”  $C^*$  such that  $C^* \cong_\varepsilon \tilde{C}$  and  $(\tilde{C}(x_i^*) \neq C^*(x_i^*)) \vee (\tilde{C}(\alpha_i^*) \neq C^*(\alpha_i^*))$  for a majority of  $i$ . Therefore,  $\mathcal{D}$  outputs 1 with probability  $\delta - \text{negl}(\lambda)$ .

If the distribution of  $(\alpha^*, x^*)$  as that of  $(\alpha', x')$  then (1) either it is not the case that  $C^* \cong_{1/\sqrt{2}+\rho/2} \tilde{C}$  in which case the test fails and  $\mathcal{D}$  outputs 0 with overwhelming probability, or (2) it is the case that  $C^* \cong_{1/\sqrt{2}+\rho/2} \tilde{C}$  in which case: each  $i$  we have

$$\Pr[(\tilde{C}(x_i^*) = C^*(x_i^*)) \wedge (\tilde{C}(\alpha_i^*) = C^*(\alpha_i^*))] \geq (1/\sqrt{2} + \rho/2)^2 = 1/2 + \Omega(\rho).$$

Moreover, since the events are independent for each  $i$ , the probability that equality holds for fewer than a majority of  $i \in [\ell]$  is (by Chernoff bound)  $2^{-\Omega(\rho^2 \ell)}$  which is negligible. Therefore,  $\mathcal{D}$  outputs 1 with probability  $\text{negl}(\lambda)$ .  $\square$

*Proof of Lemma 4.12.* This proof is essentially the same as that of Lemma 4.5. We explain only the important differences. In the message-embedding construction, for marked points, the output of a marked program is  $F(K_3, x^\perp \parallel \beta) \oplus \text{msg}$ . When moving from Hybrid 1 to Hybrid 2 we remove the marked points from the challenge circuit  $\tilde{C}$  which just outputs  $F(K, x)$  on all inputs  $x$ , but we need to patch the detection key appropriately so it still detects  $\tilde{C}$  as marked and outputs the right message. To do that we modify the intermediate hybrids  $\text{Hyb}_1^j$  so as to use the following circuits for the challenge marked circuit  $\tilde{C}$  and the detection key component  $\tilde{D}_2$  respectively:

- $C^j[K, K_1, K_2, K_3, \text{msg}]$  described in Figure 11
- $D_2^j[K, K_1, K_2, K_3, \text{msg}]$  described in Figure 12.

**Constants:** pPRF keys  $K, K_1, K_2, K_3$ , and message  $\text{msg}$ .

**Inputs:**  $x \in \mathcal{D}$

**Code:** Parse  $x = x^\perp \parallel x^\text{R}$ . Set  $\alpha := F(K_1, x^\perp)$ ,  $\beta := F(K, \alpha)$ ,  $\hat{x}^\text{R} := F(K_2, x^\perp \parallel \beta)$ .  
If  $x^\text{R} = \hat{x}^\text{R}$  and  $x^\perp \geq j$ , then output  $F(K_3, x^\perp \parallel \beta) \oplus \text{msg}$ , else output  $F(K, x)$ .

Figure 11: Program  $C^j[K, K_1, K_2, K_3, \text{msg}]$

In Hybrid 2, the detection key  $\tilde{D}_2$  uses the circuit  $D_2^j$ . The important difference here is that program  $D_2^j[K, K_1, K_2, K_3, \text{msg}]$  used in the detection key must contain the challenge message  $\text{msg}$  embedded in the challenge program. To embed  $\text{msg}$  in the hybrid detection key, we must know the message before we give  $dk$  to the adversary. Thus, the scheme achieves selective-message security. Note that the adversary can adaptively query pPRF key  $K^{(j)}$  and message  $\text{msg}_i$  to  $\mathcal{MO}(mk, \cdot, \cdot)$ . The rest of the proof is essentially the same where all the hybrids are defined analogously.  $\square$

**Constants:** pPRF keys  $K, K_1, K_2, K_3$ , and message  $\text{msg}$ .

**Inputs:**  $r, \beta$

**Code:** Set  $x^L := G(r)$ ,  $\alpha := F(K_1, x^L)$ ,  $x^R := F(K_2, x^L \parallel \beta)$ .

If  $x^L \geq j$  or  $\beta \neq F(K, \alpha)$ , then set  $y := F(K_3, x^L \parallel \beta)$ , else set  $y := F(K, x^L \parallel x^R) \oplus \text{msg}$ .

Output  $(x^R, y)$ .

Figure 12: Program  $D_2^j[K, K_1, K_2, K_3, \text{msg}]$

**Full Security.** We can achieve full security with an adaptively (rather than selectively) chosen embedded message by using the standard complexity leveraging technique where we guess the message  $\text{msg}$  before starting the security game. This loses a factor of  $|\mathcal{M}|$  in the security reduction which is super-polynomial if the message space  $\mathcal{M}$  is. However, since we already rely on complexity leveraging and assume sub-exponentially secure iO and PRFs to achieve public-detection, this does not impact our full result in any significant way.

**Note on the approximation factor.** Our scheme only achieves security for approximation factors  $\varepsilon > 1/\sqrt{2}$  rather than  $\varepsilon > 1/2$  (which we showed optimal in Section 3). This comes from the fact the detection procedure calls the circuit  $C'$  twice and requires both calls to return the correct value a majority of the time. In fact, there is a concrete attack on our message-embedding construction for an approximation factor  $\varepsilon$  in the range  $\frac{1}{2} \leq \varepsilon \leq \frac{1}{\sqrt{2}}$ . We leave it as an open problem if it is possible to construct a scheme that works for all  $\varepsilon > 1/2$ .

**Improved Rate with ECC.** We note that we can think of the above construction as relying on a naive *repetition code* over the alphabet  $\Sigma = \{0, 1\}^m$ , where we just repeat the same message many times. We can modify the construction to allow for a larger message domain  $\mathcal{M}$  than just  $\Sigma$  by using better codes, such as the Reed-Solomon code. In this variant we would set  $\mathcal{M} = \Sigma^d$  for some polynomial  $d$  (without increasing the parameters  $n, m$  of the underlying pPRF). We would think of a message  $\text{msg} \in \Sigma^d$  as a degree  $d - 1$  polynomial  $p_{\text{msg}}$  over  $\mathbb{F}_{2^m}$ . For simplicity, let's assume that  $n = m$  and in this case we can modify the program in Figure. 10 so that, instead of outputting  $F(K_3, x^L \parallel \beta) \oplus \text{msg}$  we would output  $F(K_3, x^L \parallel \beta) \oplus p_{\text{msg}}(x)$ . The detection procedure would apply the error-decoding procedure of the Reed-Solomon code to recover the message.

**List Decoding.** We note that our construction could also be modified to work for any  $\varepsilon = 1/\text{poly}(\lambda)$  approximation factors by (necessarily) relaxing the definition and allowing the detection procedure to output a (small) list of possible messages rather than a single message. For security, we would require that the correct message is contained in the list. For example, in our construction, instead of outputting the “majority value”  $\text{msg}$  such that  $|\{i : \text{msg} = \text{msg}_i\}| > \ell/2$  we could output all  $O(1/\varepsilon)$  values  $\text{msg}$  such that  $|\{i : \text{msg} = \text{msg}_i\}| > \varepsilon\ell$  (or output unmarked if no such value exists). By signing the messages with a standard signature scheme, we could also ensure that the list of messages output by the detection procedure can only contain (in addition to the correct message) the messages that were embedded in some watermarked circuit by a previous call to the marking oracle.

## 5 Watermarking Other Cryptographic Primitives

We show how to use our watermarking scheme for pPRF to watermark other cryptographic primitives. Sahai and Waters proposed a public-key encryption (PKE) scheme and a signature (SIG) scheme based on pPRF and iO [SW14]. We briefly review their schemes since we apply our watermarking schemes to their PKE and SIG schemes.

**The Sahai-Waters PKE scheme.** Let  $G$  be a PRG.

**Key Generation:** Choose pPRF key  $K \leftarrow \text{Key}(1^\lambda)$  and generate an encryption program  $\mathcal{E}$  in Figure 13.  
Output  $(pk, dk) := (i\mathcal{O}(\mathcal{E}), K)$ .

**Encryption:** For plaintext  $m$ , choose randomness  $r \in \{0, 1\}^\lambda$  and output  $ct := i\mathcal{O}(\mathcal{E})(m, r)$ .

**Decryption:** Parse  $ct = (c_1, c_2)$  and output  $m' := F(K, c_1) \oplus c_2$ .

**Constants:** pPRF key  $K$ .

**Inputs:** plaintext  $m \in \{0, 1\}^m$ , randomness  $r \in \{0, 1\}^\lambda$

**Code:** Let  $t := G(r)$  and output  $ct := (t, F(K, t) \oplus m)$ .

Figure 13: Encryption Program  $\mathcal{E}$

**The Sahai-Waters SIG scheme.** Let  $f$  be a one-way function.

**Key Generation:** Choose pPRF key  $K \leftarrow \text{Key}(1^\lambda)$  and generate a verification program  $\mathcal{V}$  in Figure 14.  
Output  $(vk, sk) := (i\mathcal{O}(\mathcal{V}), K)$ .

**Signing:** For message  $m$ , output  $\sigma := F(K, m)$ .

**Verification:** Output what  $i\mathcal{O}(\mathcal{V})(m, \sigma)$  outputs.

**Constants:** pPRF key  $K$ .

**Inputs:** Signature  $\sigma$ , message  $m$

**Code:** Check if  $f(\sigma) = f(F(K, m))$ . If it holds, then output 1, otherwise 0.

Figure 14: Verification Program  $\mathcal{V}$

As we reviewed, the decryption procedure of the PKE scheme and the signing procedure of the SIG scheme are just evaluations of a pPRF. Thus, we can watermark these primitives by using our watermarking schemes for pPRF. However, there is a subtle issue when we prove the non-removability of a marked decryption/signing program. In these settings, the adversary is given public values related to a pPRF key, that is,  $i\mathcal{O}(\mathcal{E})$  or  $i\mathcal{O}(\mathcal{V})$ . We should show the PKE and SIG schemes are secure even if  $i\mathcal{O}(\mathcal{E})$  and  $i\mathcal{O}(\mathcal{V})$  are constructed from marked pPRF key  $\tilde{K} = \text{Mark}(mk, K, \text{msg})$  to reduce the non-removability of marked decryption/signing program to non-removability of pPRF. To achieve this and for modular analysis, we introduce an extension of pPRF, *puncturable marked PRF*.

**Definition 5.1 (Puncturable Marked PRF)** A puncturable marked PRF (pmPRF) scheme from pPRF ( $\text{Key}, F, \text{Puncture}$ ) consists of watermarking scheme for pPRF ( $\text{Gen}, \text{Mark}, \text{Detect}$ ), a puncture-with-mark algorithm  $\text{PunctureM}$ , and an evaluation-with-mark algorithm  $\tilde{F}$ . They satisfy the following conditions. (we omit  $\mathcal{F}$  in subscripts in algorithms.)

**Functionality preserving under puncturing:** For all polynomial size set  $S \subseteq \{0, 1\}^{n(\lambda)}$ , for all  $x \in \{0, 1\}^{n(\lambda)} \setminus S$ , for all  $\text{msg} \in \mathcal{M}$ , it holds that

$$\Pr \left[ \begin{array}{l} \tilde{F}(\tilde{K}, x) = \tilde{F}(\tilde{K}\{S\}, x) : \\ K \leftarrow \text{Key}(1^\lambda), (dk, mk) \leftarrow \text{Gen}(1^\lambda), \\ \tilde{K} := \text{Mark}(mk, K, \text{msg}), \\ \tilde{K}\{S\} := \text{PunctureM}(mk, K, \text{msg}, S) \end{array} \right] = 1.$$

**Pseudorandom at punctured points:** For all polynomial size set  $S \subseteq \{0, 1\}^{n(\lambda)}$  and  $\text{msg} \in \mathcal{M}$ , it holds that for all PPT adversary  $\mathcal{A}$ ,

$$\left| \Pr[\mathcal{A}(\tilde{K}\{S\}, \tilde{F}(\tilde{K}, S)) = 1] - \Pr[\mathcal{A}(\tilde{K}\{S\}, U_{m(\lambda) \cdot |S|}) = 1] \right| \leq \text{negl}(\lambda)$$

where  $(dk, mk) \leftarrow \text{Gen}(1^\lambda)$ ,  $K \leftarrow \text{Key}(1^\lambda)$ ,  $\tilde{K} := \text{Mark}(mk, K, \text{msg})$ ,  $S = \{x_1, \dots, x_k\}$ ,  $\tilde{K}\{S\} := \text{PunctureM}(mk, K, \text{msg}, S)$ ,  $\tilde{F}(\tilde{K}, S)$  denotes multiple values  $(\tilde{F}(\tilde{K}, x_1), \dots, \tilde{F}(\tilde{K}, x_k))$ , and  $U_\ell$  denotes the uniform distribution over  $\ell$  bits.

We can easily construct pmPRF from our watermarking schemes by replacing pPRF key  $K$  in a marked program with a punctured key  $\tilde{K}\{S\}$ . For example,  $\text{PunctureM}(mk, K, \text{marked}, S)$  can compute a punctured and marked pPRF key  $\tilde{K}\{S\} = C[K\{S\}, \mathbf{x}, \mathbf{y}]$  in the basic construction if we have  $mk = wk = (\mathbf{x}, \mathbf{y})$  and  $K$ . Although we cannot compute  $\tilde{K}\{S\}$  from  $K\{S\}$  even if we use  $mk$ , it does not cause any problem since we use  $\tilde{K}\{S\}$  only in security proofs. Similarly, we can construct pmPRF from the constructions with public-key detection in Section 4.2.

**Security of Sahai-Waters PKE and SIG based on pmPRF.** Even if we replace pPRF  $F(K, \cdot)$  in Sahai-Waters PKE and SIG with pmPRF  $\tilde{F}(\tilde{K}, \cdot)$ , they remain secure since marked pPRFs have statistical correctness and we can puncture marked pPRF keys as we observed. Thus, the original proofs by Sahai and Waters still work.

**Constants:** Marked pPRF key  $\tilde{K}$ .

**Inputs:**  $ct = (c_1, c_2)$

**Code:** Output  $\tilde{F}(\tilde{K}, c_1) \oplus c_2$ .

Figure 15: Marked Decryption Program  $\tilde{\mathcal{D}}$

**Constants:** Marked pPRF key  $\tilde{K}$ .

**Inputs:**  $m$

**Code:** Output  $\tilde{F}(\tilde{K}, m)$ .

Figure 16: Marked Signing Program  $\tilde{\mathcal{S}}$

**On Non-Removability of Marked Decryption and Signing Programs.** If there exists adversary  $\mathcal{A}$  of the non-removability security game for PKE, then we can construct adversary  $\mathcal{B}$  of the non-removability security game for pPRF. We focus on the PKE case and the SIG case is similarly proved.  $\mathcal{B}$  is given  $\tilde{K} = \text{Mark}(mk, K, \text{msg})$ .  $\mathcal{B}$  gives  $\mathcal{A}$  program  $\tilde{\mathcal{D}}$  described in Figure 15 and a public key  $i\mathcal{O}(\tilde{\mathcal{E}})$  where  $\tilde{\mathcal{E}}$  is the same as  $\mathcal{E}$  except that  $K$  is replaced with  $\tilde{K} = \text{Mark}(mk, K, \text{msg})$ . Giving  $i\mathcal{O}(\tilde{\mathcal{E}})$  does not give the adversary any more power to break the non-removability since it is constructed from the marked pPRF key.

Eventually,  $\mathcal{A}$  outputs un-marked program  $C^*$  such that  $C^* \cong_\varepsilon \tilde{\mathcal{D}}$ . Thus,  $\mathcal{B}$  outputs a program that computes  $C^*(c_1, c_2) \oplus c_2$  for input  $(c_1, c_2)$  as an un-marked pPRF program. This is a valid attack since we assume  $\mathcal{A}$  breaks the non-removability of the marked decryption program  $\tilde{\mathcal{D}}$ .

## 6 Conclusions

We have shown how to watermark various cryptographic programs such as PRFs, public-key encryption and signatures. Our watermarking schemes have a public-key detection procedure and achieve security in the presence of the marking oracle. The attacker cannot remove the mark by coming up with any program that functionally approximates the marked program where we can achieve an approximation factor of  $\varepsilon = 1/\text{poly}$  in the case of “mark-embedding” schemes and  $\varepsilon = \frac{1}{\sqrt{2}} + 1/\text{poly}$  in the case of “message-embedding” schemes. Several interesting open questions remain.

Firstly, it would be interesting to further explore the connection between obfuscation and watermarking and to see if some form of obfuscation is *necessary* to achieve watermarking or if one can come up with constructions that avoid obfuscation. Less ambitiously, one could hope to at least avoid complexity leveraging in the constructions of watermarking.

Secondly, it would be interesting to achieve a fully public-key watermarking construction where both the marking and the detection procedure only use public keys. This kind of watermarking appears plausible and we can even come up with a heuristic construction by taking our watermarking scheme with a secret-key marking procedure and releasing an obfuscated marking procedure as a public marking key. However, proving such a scheme secure by only relying on iO appears difficult.

Thirdly, it would be interesting to achieve a message-embedding construction with the optimal approximation factor  $\varepsilon > 1/2$  rather than  $\varepsilon > 1/\sqrt{2}$ .

Lastly, it would be interesting to come up with watermarking schemes for richer classes of programs as well as more applications for watermarking cryptographic programs.

## References

- [BGI<sup>+</sup>01] Boaz Barak, Oded Goldreich, Russell Impagliazzo, Steven Rudich, Amit Sahai, Salil P. Vadhan, and Ke Yang. On the (im)possibility of obfuscating programs. In *Advances in Cryptology - CRYPTO 2001, 21st Annual International Cryptology Conference, Santa Barbara, California, USA, August 19-23, 2001, Proceedings*, pages 1–18, 2001.
- [BGI<sup>+</sup>12] Boaz Barak, Oded Goldreich, Russell Impagliazzo, Steven Rudich, Amit Sahai, Salil P. Vadhan, and Ke Yang. On the (im)possibility of obfuscating programs. *J. ACM*, 59(2):6, 2012.
- [BGI14] Elette Boyle, Shafi Goldwasser, and Ioana Ivan. Functional signatures and pseudorandom functions. In *Public-Key Cryptography - PKC 2014 - 17th International Conference on Practice and Theory in Public-Key Cryptography, Buenos Aires, Argentina, March 26-28, 2014. Proceedings*, pages 501–519, 2014. Full version available from <http://eprint.iacr.org/2013/401>.

- [BPR12] Abhishek Banerjee, Chris Peikert, and Alon Rosen. Pseudorandom functions and lattices. In *Advances in Cryptology - EUROCRYPT 2012 - 31st Annual International Conference on the Theory and Applications of Cryptographic Techniques, Cambridge, UK, April 15-19, 2012. Proceedings*, pages 719–737, 2012.
- [BW13] Dan Boneh and Brent Waters. Constrained pseudorandom functions and their applications. In *Advances in Cryptology - ASIACRYPT 2013 - 19th International Conference on the Theory and Application of Cryptology and Information Security, Bengaluru, India, December 1-5, 2013, Proceedings, Part II*, pages 280–300, 2013. Full version available from <http://eprint.iacr.org/2013/352>.
- [GGH<sup>+</sup>13] Sanjam Garg, Craig Gentry, Shai Halevi, Mariana Raykova, Amit Sahai, and Brent Waters. Candidate indistinguishability obfuscation and functional encryption for all circuits. In *54th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2013, 26-29 October, 2013, Berkeley, CA, USA*, pages 40–49, 2013. Full version available from <http://eprint.iacr.org/2013/451>.
- [GGM86] Oded Goldreich, Shafi Goldwasser, and Silvio Micali. How to construct random functions. *J. ACM*, 33(4):792–807, 1986.
- [HMW07] Nicholas Hopper, David Molnar, and David Wagner. From weak to strong watermarking. In *Theory of Cryptography, 4th Theory of Cryptography Conference, TCC 2007, Amsterdam, The Netherlands, February 21-24, 2007, Proceedings*, pages 362–382, 2007.
- [KPTZ13] Aggelos Kiayias, Stavros Papadopoulos, Nikos Triandopoulos, and Thomas Zacharias. Delegatable pseudorandom functions and applications. In *2013 ACM SIGSAC Conference on Computer and Communications Security, CCS'13, Berlin, Germany, November 4-8, 2013*, pages 669–684, 2013. Full version available from <http://eprint.iacr.org/2013/379>.
- [Nis13] Ryo Nishimaki. How to watermark cryptographic functions. In *Advances in Cryptology - EUROCRYPT 2013, 32nd Annual International Conference on the Theory and Applications of Cryptographic Techniques, Athens, Greece, May 26-30, 2013. Proceedings*, pages 111–125, 2013. Full version available from <http://eprint.iacr.org/2014/472>.
- [NSS99] David Naccache, Adi Shamir, and Julien P. Stern. How to copyright a function? In *Public Key Cryptography, Second International Workshop on Practice and Theory in Public Key Cryptography, PKC '99, Kamakura, Japan, March 1-3, 1999, Proceedings*, pages 188–196, 1999.
- [SW14] Amit Sahai and Brent Waters. How to use indistinguishability obfuscation: deniable encryption, and more. In *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 - June 03, 2014*, pages 475–484, 2014. Full version available from <http://eprint.iacr.org/2013/454>.
- [YF11] Maki Yoshida and Toru Fujiwara. Toward digital watermarking for cryptographic data. *IEICE Transactions*, 94-A(1):270–272, 2011.