

Publicly Verifiable Software Watermarking

Aloni Cohen* Justin Holmgren† Vinod Vaikuntanathan‡

December 7, 2015

Abstract

Software Watermarking is the process of transforming a program into a functionally equivalent “marked” program in such a way that it is computationally hard to remove the mark without destroying functionality. Barak, Goldreich, Impagliazzo, Rudich, Sahai, Vadhan and Yang (CRYPTO 2001) defined software watermarking and showed that the existence of indistinguishability obfuscation implies that software watermarking is impossible. Given the recent candidate constructions of indistinguishability obfuscation, this result paints a bleak picture for the possibility of meaningful watermarking.

We show that slightly relaxing the functionality requirement gives us strong positive results for watermarking. Namely, instead of requiring the marked program to agree with the original unmarked program on *all inputs*, we require only that they agree on a large fraction of inputs. With this relaxation in mind, our contributions are as follows.

1. We define publicly verifiable watermarking where marking a program requires a secret key, but anyone can verify that a program is marked. The handful of existing watermarking schemes are secretly verifiable, and moreover, satisfy only a weak definition where the adversary is restricted in the type of unmarked programs it is allowed to produce (Naccache, Shamir and Stern, PKC 1999; Nishimaki, EUROCRYPT 2013). Moreover, our definition requires security against chosen program attacks, where an adversary has access to an oracle that marks programs of her choice.
2. We construct a publicly verifiable watermarking scheme for any family of puncturable pseudo-random functions (PPRF), assuming indistinguishability obfuscation and injective one-way functions.

Complementing our positive result, we show that there are pseudo-random functions which cannot be watermarked, even in a very weak setting. As a corollary, we demonstrate the first family of PRFs that are not point-puncturable.¹

*E-mail: aloni@mit.edu. MIT.

†E-mail: holmgren@mit.edu. MIT.

‡E-mail: vinodv@mit.edu. MIT. Research supported in part by DARPA Grant number FA8750-11-2-0225, an Alfred P. Sloan Research Fellowship, Microsoft Faculty Fellowship, and a Steven and Renee Finn Career Development Chair from MIT.

¹This work has since been improved, revised, and merged with [NW15]. The merged version is [CHN⁺15].

Contents

1	Introduction	3
1.1	Our Results	3
1.2	Our Techniques	4
1.3	Related Work	5
2	Preliminaries and Definitions	6
2.1	Watermarking Schemes	6
2.1.1	Weakened Definitions	8
2.2	Remarks on the Definition	9
3	Amplifying Unremovability and Unforgeability	9
4	A Distinguishing to Removing Reduction	11
5	Main Construction	12
6	The Limits of Watermarking	14
A	Puncturable Encryption	18
A.1	Required Properties	19
B	Puncturable Encryption Construction	20
B.1	Construction	20
B.2	Ciphertext Pseudorandomness	22
C	Proof of Theorem 5.3	26
D	Proof of Theorem 5.4	29
E	Proof of Theorem 4.2	31
F	Unwatermarkable PRFs	31
F.1	Preliminaries	32
F.2	Construction	33
F.3	Learnability	34
F.4	Pseudorandomness	35
G	Relationship between γ and δ	37
G.1	Multi-bit Equivalence	38

1 Introduction

Software watermarking is the process of embedding a “mark” in a program so that the marked program preserves functionality, and furthermore, it is impossible to remove the mark without destroying functionality. Despite its numerous applications in digital rights management and copy protection, rigorous definitions of software watermarking and mathematically sound constructions have been few and far between.

Software watermarking was first formally defined in the seminal work of Barak, Goldreich, Impagliazzo, Rudich, Sahai, Vadhan and Yang [BGI⁺12]. They showed that the existence of indistinguishability obfuscation (IO) implies that software watermarking *cannot* exist (for any non-trivial class of programs). Given the recent candidate constructions of IO [GGH⁺13, BR14, PST14, GLSW14], this result paints a rather dismal picture for the possibility of meaningful watermarking. Fortunately, though, this impossibility result crucially relies on the fact that the marked version of a program *computes the same function* as the original program. Indeed, they suggested that an *approximate functionality-preserving* relaxation of the definition wherein the programs agree on a large fraction (but not all) inputs might lend itself to positive results.

In this work, we pursue the notion of approximate functionality, and show constructions of watermarking for a general class of functions, namely any family of puncturable pseudo-random functions [BW13, BGI14, KPTZ13], under a strong definition of security. Our construction is *publicly verifiable*, meaning that while marking a program requires a secret (marking) key, verification is public. Moreover, our construction is secure against *chosen program attacks*, where the adversary gets access to an oracle that marks any program of its choice. Curiously enough, our construction relies on the existence of indistinguishability obfuscators.

A natural question that arises out of our positive result is whether all families of PRFs (more generally, any class of unlearnable circuits) can be watermarked. Complementing our positive result, we show that there are pseudo-random functions which cannot be watermarked. Our result relies on a generalization of the notion of robust unobfuscatable functions of Bitansky and Paneth [BP12].

We describe our results in more detail.

1.1 Our Results

Our first contribution is to define the notion of public-key watermarking, building on the beautiful work of Hopper, Molnar and Wagner [HMW07] who introduced a secret-key definition.² Roughly speaking, in a watermarking scheme, there is a marking algorithm *Mark* that takes as input a program P and uses the (secret) marking key MK to produce a marked program $\#P$ ³ while the verification algorithm *Verify* uses the (public) verification key VK to recognize legally marked programs. A watermarking scheme should satisfy three properties:

- *Approximately Functionality Preserving*: The program $\#P$ should agree with P in at least $1 - \rho(n)$ fraction of inputs, where ρ is the approximate functionality parameter.
- *Unremovability*: We say that a watermark remover \mathcal{R} succeeds given $\#P$ if she produces as program \hat{P} that is approximately (functionally) equivalent to $\#P$ and yet, *Verify* fails

²While the work of [HMW07] targeted the watermarking of perceptual objects such as text, audio and video, the definition remains essentially unchanged for watermarking programs or circuits.

³Marked programs are always denoted as $\#P$ or $\#C$ (in the case of circuits) in this paper.

to recognize \hat{P} as a marked program. Unremovability says that this should be hard for polynomial-time removers \mathcal{R} .

- *Unforgeability*: The other side of the coin is unforgeability which requires that a forger \mathcal{F} cannot succeed in producing a new marked program, given only $\#P$ and the verification key VK .

Moreover, we require security against *chosen program attacks*, namely that unremovability and unforgeability hold even given access to a `Mark` oracle that produces marks programs for the adversary. (In this case, the forger needs to produce a marked program that is sufficiently different from all his queries to the `Mark` oracle).

Armed with this definition, our main result is the construction of a publicly verifiable watermarking scheme for any family of puncturable pseudo-random functions. Puncturable pseudo-random functions (PPRFs) [BW13, BGI14, KPTZ13] are pseudorandom functions wherein the owner of the key K can produce a punctured key K_x that allows computation of the PRF on all inputs $y \neq x$. Moreover, given the punctured key, $\text{PRF}_K(x)$ is pseudorandom. PPRFs have seen a large number of uses recently in applications of indistinguishability obfuscation. We show:

Theorem 1.1 (Informal). *Assuming indistinguishability obfuscation and injective one-way functions, there is a watermarking scheme for any family of puncturable pseudo-random functions.*

1.2 Our Techniques

We describe our publicly verifiable watermarking construction through a sequence of ideas. As a starting point, we construct a simple secret-key watermarking scheme for puncturable PRFs. Intuitively, the (secret) marking key specifies a random pair (x^*, y^*) , and the watermarking of a PRF key K is simply the indistinguishability obfuscation of a program P_{K,x^*,y^*} that does the following:

On input x , output y^ if $x = x^*$, and $\text{PRF}_K(x)$ otherwise.*

(Secret-key) verification of a circuit C is simply checking whether $C(x^*) \stackrel{?}{=} y^*$. Intuitively, it should be hard to remove the mark since the obfuscation $\#P = \mathcal{O}(P_{K,x^*,y^*})$ “should hide” the location x^* where the functionality is altered. In other words, removing the mark necessitates altering the program at essentially all points. We formalize this intuition via a two-step proof.

First, we show a general *distinguishing-to-removing* reduction. Consider a watermark remover that takes the program $\#P$ and outputs an unmarked, yet approximately equivalent, program Q . We claim that the remover can be used to distinguish between the special input x^* and a uniformly random input x . This is because (a) we know that $\#P(x^*) \neq Q(x^*)$ (simply because the watermark remover succeeded), and yet: (b) for a uniformly random input x , $\#P(x) = Q(x)$ w.h.p. (because Q and $\#P$ agree on a large fraction of inputs). This idea can be formalized as a general distinguishing-to-removing reduction, as long as the watermark verifier uses the program $\#P$ as a black-box (which is true for all our constructions).

Secondly, we show that $\#P$ hides x^* , crucially relying on both the pseudorandomness of the family as well as its puncturability. First, one can indistinguishably replace $\#P$ by a program that uses the punctured PRF key K_{x^*} to compute the PRF on all inputs $x \neq x^*$; this is because of the IO security guarantee. Second, we can replace y^* by $\text{PRF}_K(x^*)$ indistinguishably; this is because of pseudo-randomness. Finally, change the program to one that simply computes the PRF on all

points, using the IO security guarantee again. At this point, the program contains no information about x^* whatsoever.

Unforgeability can be shown using similar ideas.

However, this construction falls to the following attack, against an adversary that obtains the marked version of even a single program of its choice. Consider an adversary that obtains a marked version $\#Q$ of an unmarked program Q that she chooses. She can build a distinguisher for x^* using these two programs. Indeed, if $Q(x) \neq \#Q(x)$, then x is likely x^* . In other words, she can build a program $\text{Dist}_{Q,\#Q}$ that, on input x , predicts whether $x = x^*$. With this newfound ability, the adversary can easily remove marks from programs. Indeed, given a marked program $\#P$, it builds a wrapper around P that first checks if an input x is x^* . If yes, it outputs \perp , otherwise, it computes the program correctly.

Looking back at this attack (which we call the “majority attack”), we realize that the main source of difficulty is that we reuse the trigger point x^* across all programs. Our main idea to circumvent this attack is to make the trigger point *program-dependent*. In particular, the marking algorithm probes the program on a set of pre-determined inputs to obtain x^* (That is, x^* is a function of the *values* of the program on these inputs). Then, it goes ahead and changes the output of the program on x^* . This allows us to construct a watermarking scheme secure against *lunch-time chosen program attacks*. That is, the adversary can query the Mark oracle before it obtains the challenge program, but not after.

Finally, we augment the construction to be publicly verifiable. Our security proof sketched above relied crucially on the fact that the trigger points cannot be distinguished from uniformly random points, and yet verification seems to require knowledge of these points. We resolve this apparently conundrum by first embedding an exponentially large number (yet an exponentially small fraction) of trigger points, and observing that while verification requires a program that generates a random trigger point, security only requires that a uniformly random trigger point is pseudo-random. These requirements can indeed co-exist, and in fact, we use a variant of the *hidden sparse trigger* machinery of Sahai and Waters [SW14] to achieve both effect simultaneously. We refer the reader to Section 5 for the technical details.

Limits of Watermarking. A natural question that arises out of our work is whether all classes of circuits can be watermarked (of course, with approximate preservation of functionality).⁴ We show that there are PRF families that cannot be watermarked. Our result relies on the notion of robust unobfuscatable functions of Bitansky and Paneth [BP12]. (See Section 6 for more details.)

1.3 Related Work

There has been a large body of work on watermarking in the applied research community. Notable contributions of this line of research include the discovery of *protocol attacks* such as the copy attack by Kutter, Voloshynovskiy and Herrigel [KVH00] and the ambiguity attack by Adelsback, Katzenbeisser and Veith [AKV03]. However, these works do not formally define the security guarantees required of watermarking, and have resulted in a cat-and-mouse game of designing watermarking schemes that are broken fairly immediately.

⁴Clearly, circuits that can be exactly learned from black-box access cannot be watermarked, for obvious reasons. Thus, our question is non-trivial for unlearnable circuits.

There are a handful of works that propose rigorous notions of security for watermarking, with associated proofs of security based on complexity-theoretic assumptions. Hopper, Molnar and Wagner [HMW07] formalized strong notions of watermarking security with approximate functionality; our definitions are inspired by their work. Barak et al. [BGI⁺12] proposed simulation-based definitions of watermarking security; their main contribution is a negative result, described earlier in the introduction, which shows that indistinguishability obfuscation rules out any meaningful form of watermarking that preserves functionality exactly. The starting point of our construction is their speculative idea that relaxing this to approximate functionality might result in positive results.

In another line of work, Naccache, Shamir and Stern [NSS99] showed how to watermark a specific hash function. Nishimaki [Nis14] recently showed a similar result along these lines, watermarking a specific construction of lossy trapdoor functions (LTDF) based on bilinear groups. These works achieve a rather weak notion of watermarking. First and most important, they define a successful adversary as one that outputs a program from the class \mathcal{C} . In particular, the watermark remover for the LTDF must output an LTDF key. In contrast, we permit the watermark remover to output any (polynomial-size) circuit. Secondly, that is they are only secure as long as a bounded number of marked programs is released to the adversary. Finally, they are both *privately verifiable*.

2 Preliminaries and Definitions

Notation. We will let λ denote a security parameter throughout this paper. We will let $\{\mathbb{C}_\lambda\}_{\lambda \in \mathbb{N}}$ be the family of all circuits with domain $\{0, 1\}^{n(\lambda)}$ and range $\{0, 1\}$ for some polynomial n . We will let $\{\mathcal{C}_\lambda\}_{\lambda \in \mathbb{N}}$ be a particular family of circuits; that is $\mathcal{C}_\lambda \subseteq \mathbb{C}_\lambda$ for all $\lambda \in \mathbb{N}$.

Let $\rho : \mathbb{N} \rightarrow [0, 1]$ be a function. For circuits $C, C' \in \mathbb{C}_\lambda$, we say that $C \sim_\rho C'$ if

$$\Pr_{x \leftarrow \{0,1\}^n} [C(x) \neq C'(x)] \leq \rho(n)$$

We call such circuits “ ρ -close”. We refer to probabilistic polynomial time algorithms simply as “p.p.t. algorithms”.

2.1 Watermarking Schemes

Our definitions will generalize and refine those of Hopper, Molnar and Wagner [HMW07]. A (public-key) watermarking scheme \mathcal{W} for a family of circuits $\mathcal{C} = \{\mathbb{C}_\lambda\}_{\lambda \in \mathbb{N}}$ is a tuple of p.p.t. algorithms (Setup, Mark, Extract), where:

- $(\text{mk}, \text{vk}) \leftarrow \text{Setup}(1^\lambda)$ is a p.p.t. key generation algorithm that takes as input a security parameter $\lambda \in \mathbb{N}$ and outputs a *marking key* mk and a *verification key* vk .
- $C_\# \leftarrow \text{Mark}(\text{mk}, C)$ is a p.p.t. marking algorithm that takes as input the marking key mk and a circuit $C \in \mathbb{C}_\lambda$ and outputs a circuit $C_\#$.
- $b \leftarrow \text{Verify}(\text{vk}, C_\#)$ is a p.p.t. algorithm which takes as input the (public) verification key vk and a (possibly marked) circuit $C_\# \in \mathbb{C}$, and outputs either **accept** (1) or **reject** (0).

Note that while Mark and Verify take any circuit in \mathbb{C} as input, we only require the correctness and security properties input to Mark is from $\mathcal{C} \subseteq \mathbb{C}$.

Correctness Properties. Having defined the syntax of a watermarking scheme, we now define the desired correctness properties. First, it is functionality preserving, namely marking a circuit C does not change its functionality too much. We formalize this by requiring that the marked and unmarked circuits agree on $\rho(n)$ fraction of the domain (where n is the length of the input to the circuit). Secondly, we require that the verification algorithm always accepts a marked circuit.

Definition 2.1 (ρ -Functionality Preserving). We say that a watermarking scheme ($\text{Setup}, \text{Mark}, \text{Verify}$) is ρ -functionality preserving if for all $C \in \mathcal{C}$:

$$\text{Mark}(\text{mk}, C) \sim_{\rho} C$$

Definition 2.2 (Completeness). A watermarking scheme ($\text{Setup}, \text{Mark}, \text{Verify}$) is said to be *complete* if

$$\Pr \left[\text{Verify}(\text{vk}, \text{Mark}(\text{mk}, C)) = 1 \mid \begin{array}{l} \text{mk}, \text{vk} \leftarrow \text{Setup}(1^{\lambda}) \\ C \leftarrow \mathcal{C} \end{array} \right] \geq 1 - \text{negl}(\lambda)$$

Security Properties. We turn to the desired security properties of the watermarking scheme. We define a scheme’s “unremovability” and “unforgeability” with respect to the following security game.

The watermarking security game is defined with two helper sets \mathcal{C} and \mathfrak{M} : \mathcal{C} is the set of marked challenge programs given to a p.p.t. adversary \mathcal{A} ; \mathfrak{M} is the set of circuits given to the adversary as a response to a $\text{Mark}(\text{mk}, \cdot)$ query. We only require that the adversary cannot “unmark” a challenge program, and that the adversary cannot “forge” a mark on a program for which he has never seen a mark.

Game 2.1 (Watermarking Security). First, the challenger generates $(\text{mk}, \text{vk}) \leftarrow \text{Setup}(1^{\lambda})$ and helper sets \mathcal{C} and \mathfrak{M} initialized to \emptyset . The adversary is presented with vk and access to the following two oracles.

- A marking oracle \mathcal{O}_M which takes a circuit C and returns $\text{Mark}(\text{mk}, C)$. \mathcal{O}_M also adds C to the set \mathfrak{M} .
- A challenge oracle \mathcal{O}_C which takes no input, but samples a circuit C^* uniformly from \mathcal{C} and returns $\#C^* = \text{Mark}(\text{mk}, C^*)$. $\#C^*$ is then added to \mathcal{C} .

Finally, \mathcal{A} outputs a circuit \hat{C} .

Our ideal notion of unremovability is that no adversary can – with better than negligible probability – output a program that is δ -close to the challenge program, but on which Verify returns zero with any noticeable probability. Along the way, we will need to consider a relaxed notion: (p, δ) -unremovable. For this, we require that no adversary can – with better than negligible probability – output a program that is δ -close to the challenge program, but on which Verify returns 0 with *probability significantly greater than p* . If $p = 0$, then this coincides with the previous notion of unremovability, which we simply call δ -unremovability.

Definition 2.3 ((p, δ) -Unremovable). In the security game, we say that the adversary (p, δ) -removes if $\Pr[\text{Verify}(\text{vk}, \hat{C}) = 0] > p + \frac{1}{\text{poly}(\lambda)}$ and $\hat{C} \sim_{\delta} C'$ for some $C' \in \mathcal{C}$. \mathcal{A} ’s (p, δ) -removing advantage is the probability that \mathcal{A} (p, δ) -removes. The scheme is (p, δ) -unremovable if all p.p.t. \mathcal{A} have negligible (p, δ) -removing advantage. The scheme is δ -unremovable if it is (p, δ) -unremovable for all $p > \text{negl}(\lambda)$.

Formally, the scheme is (p, δ) -unremovable if for all p.p.t. algorithms \mathcal{A} and all polynomials $\text{poly}(\lambda)$:

$$\Pr \left[\bigwedge \left(\Pr[\text{Verify}(\text{vk}, \hat{C}) = 0] > p + \frac{1}{\text{poly}(\lambda)} \right) \mid \begin{array}{l} \text{mk}, \text{vk} \leftarrow \text{Setup}(1^\lambda), \\ \hat{C} \leftarrow \mathcal{A}^{\mathcal{O}_M, \mathcal{O}_C}(1^\lambda, \text{vk}) \end{array} \right] \leq \text{negl}(\lambda)$$

Likewise, we say a scheme is (q, γ) -unforgeable if no adversary can – with better than negligible probability – output a program that is γ -far from all marked programs previously received from the challenger, but on which `Verify` returns 1 with *probability significantly greater than q* . If $q = 0$, then this coincides with a stronger notion we call γ -unforgeability.⁵

Definition 2.4 ((q, γ) -Unforgeable). In the security game, we say that the adversary (q, γ) -forges if $\Pr[\text{Verify}(\text{vk}, \hat{C}) = 1] > q + \frac{1}{\text{poly}(\lambda)}$ and for all $C' \in \mathfrak{M} \cup \mathfrak{C}$, $\hat{C} \not\sim_\gamma C'$. \mathcal{A} 's (q, γ) -forging advantage is the probability that \mathcal{A} (q, γ) -forges. The scheme is (q, γ) -unforgeable if all p.p.t. \mathcal{A} have negligible (q, γ) -forging advantage. The scheme is γ -unforgeable if it is (q, δ) -unforgeable for all $q > \text{negl}(\lambda)$.

Formally, the scheme is (q, γ) -unforgeable if for all p.p.t algorithms \mathcal{A} for all polynomials $\text{poly}(\lambda)$:

$$\Pr \left[\bigwedge \left(\Pr[\text{Verify}(\text{vk}, \hat{C}) = 1] > q + \frac{1}{\text{poly}(\lambda)} \right) \mid \begin{array}{l} \text{mk}, \text{vk} \leftarrow \text{Setup}(1^\lambda), \\ \hat{C} \leftarrow \mathcal{A}^{\mathcal{O}_M, \mathcal{O}_C}(1^\lambda, \text{vk}) \end{array} \right] \leq \text{negl}(\lambda)$$

2.1.1 Weakened Definitions

In our main construction, we achieve a weaker version of δ -unremovability, which we call δ -lunchtime unremovability.⁶ In this weaker definition, the adversary only has access to the marking oracle before receiving the challenger, and receives only one challenger program.

Definition 2.5 ((p, δ) -Lunchtime Unremovability). δ -lunchtime unremovability is defined as above, except that we modify the security game: the adversary can query \mathcal{O}_C at most once, after which the adversary can no longer query \mathcal{O}_M .

In our main construction, we achieve a weaker notion of γ -unforgeability. In particular, we only show that a “strong” type of γ -forgery is impossible, in effect establishing only a relaxed form of unforgeability.

In order to say that \mathcal{A} γ -strong-forges, we require an additional property on \mathcal{A} 's output \hat{C} . Instead of requiring that for all marked programs $C' \in \mathfrak{M} \cup \mathfrak{C}$ received by the adversary, there is a γ fraction of the domain on which C' differs from \hat{C} , we switch the order of the quantifiers. That is, there is a γ fraction of the domain on which for all $C' \in \mathfrak{M} \cup \mathfrak{C}$, C' differs from \hat{C} .

Definition 2.6 ((q, γ) -Relaxed Unforgeability). We say that the adversary (q, γ) -strong forges if $\Pr[\text{Verify}(\text{vk}, \hat{C}) = 1] > q + \frac{1}{\text{poly}(\lambda)}$ and for all $C' \in \mathfrak{M} \cup \mathfrak{C}$,

$$\Pr[\exists C' \in \mathfrak{M} \cup \mathfrak{C} \text{ s.t. } C'(x) = \hat{C}(x) \mid x \leftarrow \{0, 1\}^n] \leq \gamma$$

The scheme is (q, γ) -relaxed unforgeable if all p.p.t. \mathcal{A} have negligible (q, γ) -strong forging advantage.

⁵This implies a sort of “meaningfulness” property ala [BGI⁺12]: that for a random circuit $C \leftarrow \mathcal{C}$, $\text{Verify}(\text{vk}, C) = 0$ with high probability.

⁶This name is in analogy to “lunchtime” chosen ciphertext attacks on encryption schemes, in which access to a decryption oracle is given only before the challenge ciphertexts.

2.2 Remarks on the Definition

Relation to definition in “From Weak to Strong Watermarking” As discussed in the introduction, [HMW07] provide a similar definition for watermarking to that provided above. One major difference is the notion of public verification – to the best of our knowledge, we are the first to put forth this notion and provide a construction in any cryptographic setting.

Another important difference is in the definitions of ρ -functionality preserving, δ -unremovability, and γ -unforgeability. The earlier definition does not consider these parameters separately: all three are conflated and represented by δ . As observed in that work, the (im)possibility of watermarking in a particular setting depend intimately on these parameters. In this work, we separate them. In doing so, we are able to achieve constant δ for negligible ρ : that is an adversary cannot remove a mark that only alters the behavior of the original negligible fraction of the domain, even by changing the functionality of the program on 1/4 of the domain! Additionally, we show that if a watermarking scheme is both $(0, \delta)$ -unremovable and $(0, \gamma)$ -unforgeable, then $\gamma \geq \delta + \frac{1}{\text{poly}(n)}$ for some polynomial poly. (see Appendix G). We also provide an Amplification Lemma (see section 3) which depends crucially on δ and γ being distinct.

Extractable Watermarking We consider a setting in which a program is either “marked” or it is “unmarked.” A natural extension is to allow a program to be “marked with a string.” That is, Mark would take a string M as an additional argument and output a marked program C_M . A corresponding detection algorithm Extract would then extract the embedded mark. Indeed, similar definitions are considered in [BGI⁺12] and [HMW07]. The two notions coincide the space of possible marks to embed is just a singleton. The extractable setting presents a number of additional challenges, and is a very interesting direction for future work.

Perfectly Functionality Preserving. As observed in Barak et al. [BGI⁺12], if indistinguishability obfuscation exists for \mathbb{C} , then no watermarking scheme can exist that perfectly preserves functionality. To see this, choose $C \leftarrow \mathbb{C}_\lambda$ and mark it to get $C_\#$. Obfuscate both circuits (appropriately padded) – yielding \tilde{C} and $\tilde{C}_\#$. By unforgeability, $\text{Verify}(\text{vk}, \tilde{C}) = \perp$; by unremovability $\text{Verify}(\text{vk}, \tilde{C}_\#) = m$. In this way, Verify gives a distinguisher for the two obfuscated circuits. Because Mark perfectly preserves functionality, this violates the security of the obfuscation.

Even more basically, if we restrict our attention to black-box verification (as discussed in section 4), then a watermarking scheme must clearly change the functionality in some way.

3 Amplifying Unremovability and Unforgeability

Ideally, we would like to construct a watermarking scheme that that is δ -unremovable and γ -unforgeable; that is, for all PPT algorithms \mathcal{A} , the probability that \mathcal{A} can remove or forge is negligible. In this section, we show that it suffices to prove something weaker. Informally, we show that given a watermarking scheme (Setup, Mark, Verify) that is (p, δ) -unremovable and (q, γ) -unforgeable for some parameters $(1 - p) \geq q + \frac{1}{\text{poly}(\lambda)}$, we construct a watermarking (Setup, Mark, Verify') that is δ -unremovable and γ -unforgeable.

To achieve this, we amplify both security guarantees by repeating Verify a polynomial number of times, and choosing an appropriate verification threshold $\tau = \frac{q + (1 - p)}{2}$: if Verify = 1 on more

than τ trials, we return 1; if $\text{Verify} = 1$ on fewer than τ trials, we return 0. More formally, we prove the following lemma.

Lemma 3.1 (Amplification Lemma). *Let $(\text{Setup}, \text{Mark}, \text{Verify})$ be a (p, δ) -unremovable and (q, γ) -unforgeable watermarking scheme, where the run-time of Verify is t_{Verify} . Suppose that $(1-p) = q + \alpha$ for some non-negligible α . Then there is a watermarking scheme $(\text{Setup}, \text{Mark}, \text{Verify}_{p,q})$ that is $(0, \delta)$ -unremovable and $(0, \gamma)$ -unforgeable, and the run-time of $\text{Verify}_{p,q}$ is $O(\frac{\lambda}{\alpha^2}) \cdot t_{\text{Verify}}$.*

Remark 1. The proof that follows is essentially an application of a Hoeffding bound. The argument also holds for the weaker definitions of unremovability and unforgeability that we consider in this work. Specifically, in the setting of , we use the Amplification Lemma to construct a δ -lunchtime unremovable, γ -relaxed unforgeable watermarking scheme from a $(2\delta, \delta)$ -lunchtime unremovable, $(1 - \gamma, \gamma)$ -relaxed unforgeable scheme.

Proof. Let $\tau = \frac{q+(1-p)}{2}$. We define $\text{Verify}_{p,q}(\text{vk}, C)$ the following algorithm that repeatedly evaluates $\text{Verify}(\text{vk}, C)$ a total of $T = \frac{8\lambda}{\alpha^2}$ times, with independent randomness on each trial.

$$\text{Verify}_{p,q}(\text{vk}, C) \begin{cases} 1 & \text{if } \frac{1}{T} \sum_{i=1}^T \text{Verify}(\text{vk}, C) > \tau \\ 0 & \text{if } \frac{1}{T} \sum_{i=1}^T \text{Verify}(\text{vk}, C) \leq \tau \end{cases}$$

As Setup and Mark are unchanged, the new scheme preserves functionality to the same extent as the original.

To show completeness, we observe that for all $C \in \mathcal{C}$ and $C_{\#} \leftarrow \text{Mark}(\text{mk}, C)$, $\text{Verify}_{p,q}(\text{vk}, C_{\#}) = 0$ if at least a $(1 - \tau)$ fraction of the independent executions of $\text{Verify}(\text{vk}, C_{\#})$ return 0. By the completeness of the original scheme, $\Pr[\text{Verify}(\text{vk}, C_{\#}) = 0] \leq \text{negl}(\lambda)$ in each independent run. Observing that $1 - \tau$ is significantly larger than $\frac{\alpha}{4}$, and applying a standard Hoeffding bound yields $\Pr[\text{Verify}_{p,q}(\text{vk}, C_{\#}) = 0] \leq e^{-\lambda}$.

Next, we must show that $(\text{Setup}, \text{Mark}, \text{Verify}_{p,q})$ is $(0, \delta)$ -unremovable. That is, for all p.p.t. \mathcal{A} playing the watermarking security game and outputting \hat{C} , if \hat{C} is δ -close to a challenge program $C' \in \mathcal{C}$, then $\Pr[\text{Verify}_{p,q}(\text{vk}, \hat{C}) = 0]$ is negligible. Equivalently, we must show that $\Pr[\sum_{i=1}^T \text{Verify}(\text{vk}, \hat{C}) \leq \tau \cdot T]$ is negligible. By the (p, δ) -unremovability of the watermarking scheme, we know that $\Pr[\text{Verify}(\text{vk}, \hat{C}) = 1] > 1 - p - \text{negl}(\lambda)$, lower-bounding the one-shot probability of verification. Combined with the definition of τ , this implies that $\Pr[\text{Verify}(\text{vk}, \hat{C}) = 1] - \tau$ is significantly larger than $\frac{\alpha}{4}$. Applying a standard Hoeffding bound yields $\Pr[\text{Verify}_{p,q}(\text{vk}, \hat{C}) = 0] \leq e^{-\lambda}$.

Lastly, we must show that $(\text{Setup}, \text{Mark}, \text{Verify}_{p,q})$ is $(0, \gamma)$ -unforgeable. That is, for all p.p.t. \mathcal{A} playing the watermarking security game and outputting \hat{C} , if \hat{C} is γ -far from all marked programs in $\mathfrak{M} \cup \mathcal{C}$, then the probability $\Pr[\text{Verify}_{p,q}(\text{vk}, \hat{C}) = 1]$ is negligible. Equivalently, we must show that the probability $\Pr[\sum_{i=1}^T \text{Verify}(\text{vk}, \hat{C}) > \tau \cdot T]$ is negligible. By the (q, γ) -unforgeability of the watermarking scheme, we know that $\Pr[\text{Verify}(\text{vk}, \hat{C}) = 1] < q + \text{negl}(\lambda)$, upper-bounding the one-shot probability of verification. Combined with the definition of τ , this implies that $\tau - \Pr[\text{Verify}(\text{vk}, \hat{C}) = 1]$ is significantly larger than $\frac{\alpha}{4}$. Applying a standard Hoeffding bound yields $\Pr[\text{Verify}_{p,q}(\text{vk}, \hat{C}) = 1] \leq e^{-\lambda}$. □

4 A Distinguishing to Removing Reduction

The way in which we prove unremovability is by showing that no adversary can distinguish points queried by `Verify` from random. This technique is clearly restricted to schemes in which there is a notion of “points queried by `Verify`”. We will therefore restrict our attention in this work to watermarking schemes with black-box verification.

Definition 4.1 (Watermarking scheme with black-box `Verify`). We say a watermarking scheme has a **black-box verification** if `Verify`(vk, P) can be efficiently evaluated with oracle access to P .

The distinguishing-to-removing reduction we now illustrate applies to any game in which the adversary’s goal is to unmark a random marked program, but for concreteness we only show a reduction to [Theorem 2.1](#).

We now give sufficient conditions for a watermarking scheme to be (p, δ) -unremovable. In particular, the condition is that no p.p.t. algorithm \mathcal{A} has non-negligible advantage in the following game.

Game 4.1 (Distinguishing). First, the challenger samples (mk, vk) from `Setup`(1^λ). The adversary is then given vk and access to the following two oracles.

- A marking oracle \mathcal{O}_M which takes a circuit C and returns `Mark`(mk, C).
- A challenge oracle \mathcal{O}_C which when queried, samples a circuit C^* uniformly from \mathcal{C} and computes $C_\#^* = \text{Mark}(\text{mk}, C^*)$. The adversary must query \mathcal{O}_C exactly once.

At the end of the game, the challenger executes `Verify`($\text{vk}, C_\#^*$). From the points at which `Verify` queried $C_\#^*$, the challenger picks x_0 uniformly at random. The challenger also chooses x_1 as a random point in the domain of $C_\#^*$. The challenger picks a random bit b and sends x_b to the adversary.

The adversary then outputs a bit b' and wins if $b' = b$.

Lemma 4.2 (Distinguishing to Removing Reduction). *If all p.p.t. algorithms \mathcal{A} have negligible advantage in [Theorem 4.1](#), and if `Verify` queries at most L points, then `(Setup, Mark, Verify)` is $(L \cdot \delta, \delta)$ -unremovable for all δ .*

We prove the lemma by assuming the existence of an $(L \cdot \delta, \delta)$ -remover, and constructing a distinguisher for the above game. Consider a removing adversary restricted to changing at most δ fraction of the marked challenge program. If he can remove the mark with too high of a probability (significantly greater than $L \cdot \delta$), then we must show that he can distinguish points queried by `Verify` from random points with non-negligible advantage. This follows because `Verify` is black-box; intuitively, the unmarked program \hat{C} will disagree with the challenge $C_\#^*$ at x_0 with significantly greater probability than on a uniformly random point x_1 .

Remark 2. It is natural to wonder whether a converse of the above lemma holds. That is, if there exists an algorithm \mathcal{A} with non-negligible advantage in [Theorem 4.1](#), does there exist an efficient (p, δ) -remover for some δ and p ? A weak converse may be shown: if \mathcal{A} distinguishes with probability $1 - \text{negl}(\lambda)$, then a $(1 - \text{negl}(\lambda), \delta)$ -remover can be constructed for some δ that depends on the specifics of the watermarking scheme.

5 Main Construction

Theorem 5.1. *There is a watermarking scheme which is both δ -lunchtime unremovable and γ -relaxed unforgeable, for any choice of δ, γ satisfying $\gamma > 2\delta + \frac{1}{\text{poly}(\lambda)}$ for some polynomial poly.*

Proof. In [Theorem 5.3](#), we show that [Theorem 5.2](#) is $(2\delta, \delta)$ -lunchtime unremovable for every δ . In [Theorem 5.4](#), we show that [Theorem 5.2](#) is $(1 - \gamma, \gamma)$ -relaxed unforgeable for every γ .

Given any δ', γ' satisfying $\gamma' > 2\delta' + \frac{1}{\text{poly}(\lambda)}$, we note that our scheme is $(2\delta', \delta')$ -lunchtime unremovable and $(1 - \gamma', \gamma')$ -relaxed unforgeable. [Theorem 3.1](#) then implies that [Theorem 5.2](#) can be amplified into a scheme which is simultaneously $(0, \delta')$ -lunchtime unremovable and $(0, \gamma')$ -relaxed unforgeable. \square

Remark 3. In [subsection G.1](#), we describe how to extend this construction to pseudorandom function families with arbitrary output bit-length, with a small loss in parameters.

Our construction marks any PRF family $\{\mathcal{P}_\lambda : \{0, 1\}^n \rightarrow \{0, 1\}^m\}_{\lambda \in \mathbb{N}}$ that is puncturable at a single point, if $n = n(\lambda)$ and $m = m(\lambda)$ are $\Omega(\lambda^\epsilon)$ for some $\epsilon > 0$. Our construction uses the following building blocks.

- Pseudorandom function families $\{\mathcal{F}_\lambda : \{0, 1\}^n \rightarrow \{0, 1\}^m\}_{\lambda \in \mathbb{N}}$ and $\{\mathcal{G}_\lambda : \{0, 1\}^\ell \rightarrow \{0, 1\}^m\}_{\lambda \in \mathbb{N}}$ which are selectively puncturable on any interval. We refer the reader to [\[BW13\]](#) for definitions and constructions of selectively secure puncturable PRF families.
- A puncturable encryption system \mathcal{PE} with plaintexts in $\{0, 1\}^\ell$ and ciphertexts in $\{0, 1\}^n$, where $\ell = \ell(\lambda)$ is $\Omega(\lambda^\epsilon)$. We explicitly denote the randomness used in $\mathcal{PE}.\text{Enc}$ as r . We construct such a \mathcal{PE} in [Appendix A](#).
- A collision resistant hash function $\{h_\lambda : \{0, 1\}^m \rightarrow \{0, 1\}^{\ell/2}\}_{\lambda \in \mathbb{N}}$.
- A pseudorandom generator $PRG_1 : \{0, 1\}^{\ell/4} \rightarrow \{0, 1\}^{\ell/2}$, as well as another pseudorandom generator $PRG_2 : \{0, 1\}^{\ell/2} \rightarrow \{0, 1\}^n$.

Construction 5.2 (Unamplified). • **Setup** (1^λ) : Setup samples $(DK, EK) \leftarrow \mathcal{PE}.\text{Gen}(1^\lambda)$ and puncturable PRFs $F \leftarrow \mathcal{F}_\lambda$ and $G \leftarrow \mathcal{G}_\lambda$. Setup then outputs (mk, vk) , where $\text{mk} = (DK, F, G)$ and vk is the iO-obfuscation of the program in [Figure 3](#).

- **Mark** (mk, C) : Mark outputs the iO obfuscation of the circuit $C_\#$, which is described in [Figure 1](#), but padded to be as big as the largest of the circuits in any of the hybrids.
- **Verify** (vk, C) : Verify samples uniformly random bit strings $a \in \{0, 1\}^{\ell/4}$ and randomly samples r . Verify then computes $b = h(C(PR_2(PR_1(a))))$. Verify evaluates vk (an obfuscated program) on $(a||b, r)$ to obtain a pair (x, y) , and returns 1 if $C(x) = y$; otherwise, it returns 0.

As Verify is an algorithm which runs an obfuscated program vk as a subroutine, we provide the “unrolled” verification algorithm in [Figure 3](#).

We now state the main theorems of this section, and defer the proofs to appendices [C](#) and [D](#).

Theorem 5.3. *[Theorem 5.2](#) is $(2\delta, \delta)$ -unremovable for every δ .*

Theorem 5.4. *[Theorem 5.2](#) is $(1 - \gamma, \gamma)$ -relaxed unforgeable for every γ .*

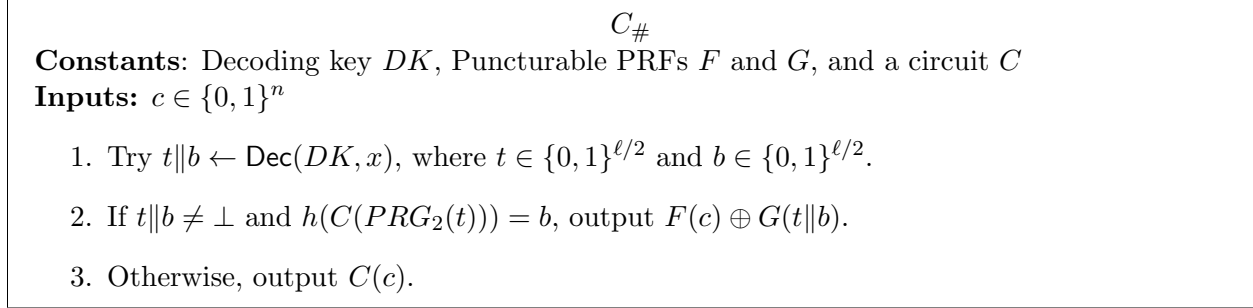


Figure 1: Program $C_{\#}$

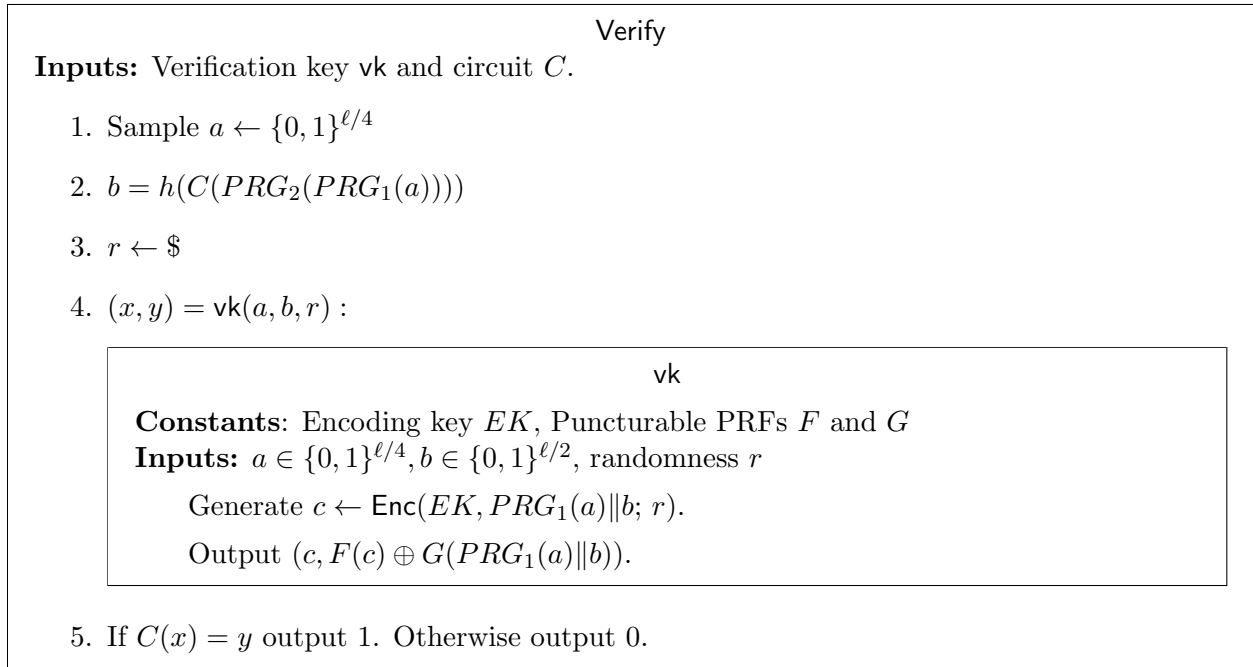


Figure 2: Unrolled $\text{Verify}(\text{vk}, C)$. Line 4 expands the execution of the program vk , which is itself an obfuscated program.

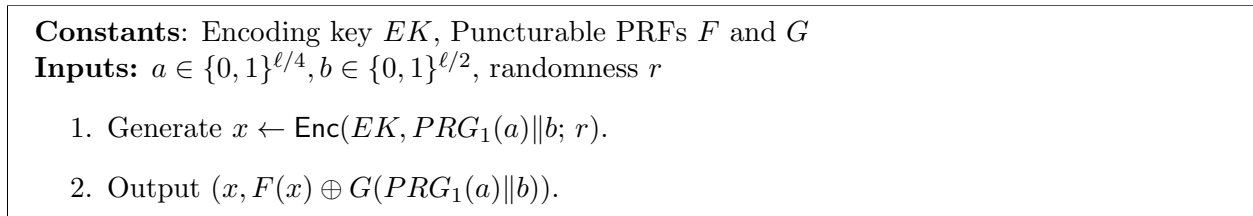


Figure 3: Verification key vk (pre-obfuscated)

6 The Limits of Watermarking

A natural question is whether there are families of functions that for which there does not exist any watermarking scheme (waterproof). Barak et al observed that general-purpose indistinguishability obfuscation rules out a notion of watermarking that *exactly* preserves functionality, but not watermarking schemes that change functionality on even a negligible fraction of the domain (as in section 5).

In this section, we discuss a number of conditions sufficient to prove that a family of circuits cannot even be watermarked – even with a much weaker form of unremovability. Informally, if a family is (non-black-box) learnable given access to a ρ -approximation of a circuit in the family, then the family is waterproof. Because it suffices to learn the family with an approximate implementation, we focus on non-black-box learnability. For such a family, from a challenge marked program the learning algorithm is able to recover the original (unmarked) program. This violates unremovability of a watermarking scheme.⁷

We construct waterproof PRFs using techniques closely related to the unobfuscatable function families of [BGI⁺12] and [BP12].⁸ In doing so, we present a construction of a family of PRFs that is not point-puncturable. To the best of our knowledge, this is the first such “unpuncturable” family of PRFs.

Consider an indexed family of functions $\mathcal{F} = \{f_K\}$ where each function is indexed by the key K . In our setting, the learning algorithm will be given any circuit g that is a $\rho(n)$ -approximate implementation of f_K , a uniformly sampled function from the family. The (randomized) learner will then output some “hypothesis” function h . If h is sufficiently close to f_K , then we can conclude that the family \mathcal{F} cannot be watermarked.

As a warm up, we begin with a very strong notion of learnability, in which the learning algorithm – here called an *extractor* – can not only output a hypothesis h which agrees with f_K on all inputs, but output the circuit f_K itself.

Definition 6.1 (ρ -Robustly Extractable Families). Let $\mathcal{F} = \{\mathcal{F}_n\}_{n \in \mathbb{N}}$ be a circuit ensemble where each family $\mathcal{F}_n = \{f_K\}_{K \in \{0,1\}^n}$. We say that \mathcal{F} is ρ -robustly extractable if there exists an efficient extractor E such that for all large enough $n \in \mathbb{N}$ and random $K \leftarrow \{0,1\}^n$, E extracts K from any circuit C such that $C \sim_{\rho(n)} f_K$:

$$\Pr[K \leftarrow E(C, 1^n)] \text{ is non-negligible.}$$

Theorem 6.1. *If \mathcal{F} is ρ -robustly extractable, then there does not exist a watermarking scheme that is $\rho(n)$ -functionality preserving, δ -unremovable, and γ -unforgeable for any $\delta \geq 0$ and $\gamma \leq 1$.⁹*

Proof. Given a ρ functionality-preserving watermarking scheme for the family \mathcal{F} , for all circuits $f_K \in \mathcal{F}$, the marked program $\#f_K = \text{Mark}(\text{mk}, f_K)$ is ρ -close to the original. If \mathcal{F} is ρ -robustly extractable, then given a challenge marked program,¹⁰ the extractor E outputs f_K with noticeable

⁷We avoid the language of learning theory because the learnability conditions we consider are not among the common settings of that field.

⁸Specifically, we extend Theorem 4.3 from [BGI⁺12] to a more general notion of approximate obfuscators.

⁹More precisely, $(0, \delta)$ -unremovability and $(1 - \frac{1}{\text{poly}(n)}, \gamma)$ -unforgeability, for any $\delta \geq 0$, $\gamma \leq 1$, and polynomial $\text{poly}(n)$.

¹⁰Recall that unremovability requires that the mark cannot be removed from a random challenge.

probability. Unless $\text{Verify}(\text{vk}, f_K) = 1$ with high probability, the extractor E violates even 0-unremovability. Otherwise we may trivially violate 1-unforgeability by simply outputting a random $f_K \leftarrow \mathcal{F}$ (without ever receiving a marked program). \square

Note that the watermarking adversary presented in the proof breaks even a very weak notion of watermarking. Namely, the adversary requires only a challenge, and no calls to a `Mark` oracle nor to a `Verify` key or verification oracle whatsoever. In this very weak setting, the construction in [Theorem 5.2](#) is secure for any family of point-puncturable PRFs.

Towards proving the main theorems of this section, we weaken [Definition 6.1](#) in two ways. Combined together in [Definition 6.2](#), these weaker notions of learnability will capture richer functionalities and allow us to construct a PRF family that cannot be watermarked ([Theorem 6.3](#)). The following discussion motivates the stronger definition and outlines the proof of the corresponding [Theorem 6.2](#)

Learnable versus extractable. What if the family is only “learnable,” but not “extractable:” instead of outputting f_K itself, the learning algorithm $L(C)$ can only output a circuit h that was functionally equivalent to f_K ? One might think that this is indeed sufficient to prove [Theorem 6.1](#), but the proof encounters a difficulty.

As before, we run the learner on the challenge program to get $h = L(\#f_K)$; if $\text{Verify}(\text{vk}, h) = 0$ with noticeable probability, then unremovability is violated. On the other hand, if $\text{Verify}(\text{vk}, h) = 1$ with high probability, how is unforgeability violated? In the extractable setting, it was possible to sample a program which verifies *without ever seeing a marked version*, simply by picking $f_K \leftarrow \mathcal{F}$. In the weaker learnable setting, we only know how to sample from this verifying distribution by evaluating $L(\#f_K)$ on a marked program. But a forger must output a marked program that is substantially different (at least γ -far) than all other marked programs seen.

To get around this issue, we consider families that are learnable with *implementation independence*; that is, for any g and g' which are both ρ approximations of $f_K \in \mathcal{F}$, the distributions $L(g)$ and $L(g')$ are computationally indistinguishable.¹¹ To complete the above proof, a forger will simply evaluate $h \leftarrow L(f_K)$ for random (unmarked) f_K (rather than on the marked $\#f_K$). Input independence of L guarantees that $\text{Verify}(\text{vk}, h) = 1$ with high probability.

Approximate versus exact learning. In [Definition 6.1](#) (and the discussion above), we required that an algorithm learning a family \mathcal{F} is able to exactly recover the functionality f_K , when given g that ρ -approximates f_K . What can we prove if $h = L(g)$ is only required to δ -approximate the original function f_K ?

Though we cannot violate 0-unremovability, we might hope to violate δ -unremovability. For ρ -functionality preserving watermarking scheme, when given a marked program $\#f_K$, the learning algorithm returns $h = L(\#f_K)$ which is a δ -approximation (with noticeable probability). By similar reasoning to [Theorem 6.1](#), it must be that either $\text{Verify}(\text{vk}, h) = 1$ with high probability or δ -unremovability is violated. If L is implementation independent as above, we contradict unforgeability.

¹¹Weaker notions likely suffice because unforgeability only requires noticeable probability of forging whereas this condition gives us high probability. We consider the input independence notion because it is a simple, natural and, as we will see, powerful case.

Definition 6.2 (ρ -Robustly, δ -Approximately Implementation Independent Learnable Families). Let $\mathcal{F} = \{\mathcal{F}_n\}_{n \in \mathbb{N}}$ be a circuit ensemble where each family $\mathcal{F}_n = \{f_K\}_{K \in \{0,1\}^n}$. We say that \mathcal{F} is ρ -robustly, δ -approximately learnable if there exists an efficient learner L such that for all large enough $n \in \mathbb{N}$, random $K \leftarrow \{0,1\}^n$, and any circuit C such that $C \sim_{\rho(n)} f_K$:

$$\Pr[h \sim_{\delta} f_K : h \leftarrow L(C, 1^n)] \text{ is non-negligible.}$$

We say that L is *implementation independent* if for all C_1, C_2 that are both $\rho(n)$ -close to f_K , the distributions $L(C_1, 1^n)$ and $L(C_2, 1^n)$ are computationally indistinguishable.

Theorem 6.2. *If \mathcal{F} is ρ -robustly, δ -approximately learnable with implementation independence, then there does not exist a watermarking scheme (Setup, Mark, Verify) that is $\rho(n)$ -functionality preserving, δ -unremovable, and γ -unforgeable for any $\gamma \leq 1$.¹²*

The existence of indistinguishability obfuscation implies a 0-robust, exact, implementation independent learning algorithm for polynomial-sized circuits.¹³ Therefore the theorem rules out exact watermarking schemes, assuming the existence of iO – capturing the impossibility of exact watermarking originally presented in [BGI⁺12].

Already, this rules out watermarking a large array of families. For instance, any family that is improperly PAC learnable cannot be watermarked for any negligible function ρ . The main result of this section is that there exists a PRF family that is learnable as in Definition 6.2; the construction and proof are presented in Appendix F.

Theorem 6.3. *Assuming one-way functions, there exists a pseudorandom function family \mathcal{F}_{δ} that is ρ -robustly, δ -approximately learnable with implementation independence, for any $\delta = \frac{1}{\text{poly}(n)}$, and any negligible $\rho(n)$.*

Corollary 6.4. *Assuming one-way functions, for any negligible function $\rho(n)$, inverse-polynomial function $d(n)$, and $\gamma \leq 1$, there is a family of pseudorandom functions that is not $(\rho, \frac{1}{d(n)}, \gamma)$ -watermarkable.*

As discussed, the watermarking adversary in Theorem 6.2 breaks even a very weak notion of watermarking. The adversary requires only a challenge, and no calls to a Mark oracle nor to a Verify key or verify oracle whatsoever. In this very weak setting, the construction in Theorem 5.2 is secure for any family of point-puncturable PRFs. Thus we construct an “unpuncturable” family.

Corollary 6.5. *Assuming one-way functions, there exists a family of pseudorandom functions that is not puncturable at points.¹⁴*

¹²More precisely, $(0, \delta)$ -unremovability and $(1 - \frac{1}{\text{poly}(n)}, \gamma)$ -unforgeability, for any $\delta \geq 0$, $\gamma \leq 1$, and polynomial $\text{poly}(n)$.

¹³Observed by Nir Bitansky.

¹⁴In fact, a much simpler family already achieves this notion. Only a slight modification of the families presented in [BGI⁺12] is needed. Furthermore, we actually show that there exists a family of pseudorandom functions that is not puncturable on negligible-sized sets.

References

- [AKV03] André Adelsbach, Stefan Katzenbeisser, and Helmut Veith. Watermarking schemes provably secure against copy and ambiguity attacks. In Moti Yung, editor, *Proceedings of the 2003 ACM workshop on Digital rights management 2003, Washington, DC, USA, October 27, 2003*, pages 111–119. ACM, 2003.
- [BGI⁺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.
- [BP12] Nir Bitansky and Omer Paneth. From the impossibility of obfuscation to a new non-black-box simulation technique. In *53rd Annual IEEE Symposium on Foundations of Computer Science, FOCS 2012, New Brunswick, NJ, USA, October 20-23, 2012*, pages 223–232. IEEE Computer Society, 2012.
- [BR14] Zvika Brakerski and Guy N. Rothblum. Virtual black-box obfuscation for all circuits via generic graded encoding. In Yehuda Lindell, editor, *Theory of Cryptography - 11th Theory of Cryptography Conference, TCC 2014, San Diego, CA, USA, February 24-26, 2014. Proceedings*, volume 8349 of *Lecture Notes in Computer Science*, pages 1–25. Springer, 2014.
- [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.
- [CHN⁺15] Aloni Cohen, Justin Holmgren, Ryo Nishimaki, Vinod Vaikuntanathan, and Daniel Wichs. Watermarking cryptographic capabilities. Cryptology ePrint Archive, Report 2015/1096, 2015. <http://eprint.iacr.org/>.
- [GGH⁺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. IEEE Computer Society, 2013.
- [GLSW14] Craig Gentry, Allison B. Lewko, Amit Sahai, and Brent Waters. Indistinguishability obfuscation from the multilinear subgroup elimination assumption. *IACR Cryptology ePrint Archive*, 2014:309, 2014.
- [HMW07] Nicholas Hopper, David Molnar, and David Wagner. From weak to strong watermarking. In Salil P. Vadhan, editor, *Theory of Cryptography, 4th Theory of Cryptography Conference, TCC 2007, Amsterdam, The Netherlands, February 21-24, 2007, Proceedings*, volume 4392 of *Lecture Notes in Computer Science*, pages 362–382. Springer, 2007.

- [KPTZ13] Aggelos Kiayias, Stavros Papadopoulos, Nikos Triandopoulos, and Thomas Zacharias. Delegatable pseudorandom functions and applications. In Ahmad-Reza Sadeghi, Virgil D. Gligor, and Moti Yung, editors, *2013 ACM SIGSAC Conference on Computer and Communications Security, CCS'13, Berlin, Germany, November 4-8, 2013*, pages 669–684. ACM, 2013.
- [KVH00] M. Kutter, S. Voloshynovskiy, and A. Herrigel. The watermark copy attack. In *Proceedings of the SPIE, Security and Watermarking of Multimedia Contents II*, volume 3971, pages 371–379, 2000.
- [Nis14] Ryo Nishimaki. How to watermark cryptographic functions. *IACR Cryptology ePrint Archive*, 2014:472, 2014.
- [NSS99] David Naccache, Adi Shamir, and Julien P. Stern. How to copyright a function? In Hideki Imai and Yuliang Zheng, editors, *Public Key Cryptography, Second International Workshop on Practice and Theory in Public Key Cryptography, PKC '99, Kamakura, Japan, March 1-3, 1999, Proceedings*, volume 1560 of *Lecture Notes in Computer Science*, pages 188–196. Springer, 1999.
- [NW15] Ryo Nishimaki and Daniel Wichs. Watermarking cryptographic programs against arbitrary removal strategies. *Cryptology ePrint Archive*, Report 2015/344, 2015. <http://eprint.iacr.org/>.
- [PST14] Rafael Pass, Karn Seth, and Sidharth Telang. Indistinguishability obfuscation from semantically-secure multilinear encodings. In Juan A. Garay and Rosario Gennaro, editors, *Advances in Cryptology - CRYPTO 2014 - 34th Annual Cryptology Conference, Santa Barbara, CA, USA, August 17-21, 2014, Proceedings, Part I*, volume 8616 of *Lecture Notes in Computer Science*, pages 500–517. Springer, 2014.
- [SW14] Amit Sahai and Brent Waters. How to use indistinguishability obfuscation: deniable encryption, and more. In David B. Shmoys, editor, *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 - June 03, 2014*, pages 475–484. ACM, 2014.

A Puncturable Encryption

One of our main abstractions is a *puncturable encryption* system. This is a public-key encryption system in which the decryption key can be punctured on a set of ciphertexts. We will rely on a strong ciphertext pseudorandomness property which holds even given access to a punctured decryption key. We will additionally require that valid ciphertexts are *sparse*, and that a decryption key punctured at some set of ciphertexts C is functionally equivalent to the nonpunctured decryption key, except possibly on C .

In this section we define the puncturable encryption abstraction that we use in [section 5](#). We instantiate this definition in [Theorem A.1](#).

Syntactically, a puncturable encryption scheme \mathcal{PE} for a message space $\mathcal{M} = \{0, 1\}^\ell$ is a triple of probabilistic algorithms ($\text{Gen}, \text{Puncture}, \text{Enc}$) and a deterministic algorithm Dec . The space of ciphertexts will be $\{0, 1\}^n$ where $n = \text{poly}(\ell, \lambda)$. For clarity and simplicity, we will restrict our

exposition to the case when $\lambda = \ell$. The puncturable encryption scheme has an additional parameter $Q \in \mathbb{N}$ which determines the number of points that can be punctured.¹⁵

- $\text{Gen}(1^\lambda) \rightarrow EK, DK$: Gen takes the security parameter in unary, and outputs an encryption key EK and a decryption key DK .
- $\text{Puncture}(DK, C) \rightarrow DK\{C\}$: Puncture takes a decryption key DK , and a set of ciphertexts $C \subset \{0, 1\}^n$, of size at most Q .¹⁶ Puncture outputs a “punctured” decryption key $DK\{C\}$.
- $\text{Enc}(EK, m) \rightarrow c$: Enc takes an encryption key EK and a message $m \in \{0, 1\}^\ell$, and outputs a ciphertext c in $\{0, 1\}^n$.
- $\text{Dec}(DK, c) \rightarrow m$ or \perp : Dec takes a possibly punctured decryption key DK and a string $c \in \{0, 1\}^n$. It outputs a message m or \perp .

A.1 Required Properties

Correctness We require that for all messages m ,

$$\Pr \left[\text{Dec}(DK, c) = m \mid \begin{array}{l} (EK, DK) \leftarrow \text{Gen}(1^\lambda), \\ c \leftarrow \text{Enc}(EK, m) \end{array} \right] = 1$$

Punctured Correctness We also require¹⁷ the same to hold for keys which are punctured. For all possible keys $(EK, DK) \leftarrow \text{Gen}(1^\lambda)$, all sets $C \subset \{0, 1\}^n$ of size at most Q , all punctured keys $DK' \leftarrow \text{Puncture}(DK, C)$, and all potential ciphertexts $c \in \{0, 1\}^n \setminus C$:

$$\text{Dec}(DK, c) = \text{Dec}(DK', c)$$

Ciphertext Pseudorandomness We require that in the following game, all PPT adversaries \mathcal{A} have negligible advantage.

Game A.1 (Ciphertext Pseudorandomness).

1. \mathcal{A} sends a messages $m_1, \dots, m_{Q/2} \in \mathcal{M}$ to the challenger.
2. The challenger does the following:
 - Samples $(EK, DK) \leftarrow \text{Gen}(1^\lambda)$
 - Computes encryptions $c_i \leftarrow \text{Enc}(EK, m_i)$ for each $i \in [Q/2]$. Let $\vec{c} = (c_1, \dots, c_{Q/2})$.
 - Samples $r_1, \dots, r_{Q/2} \leftarrow \{0, 1\}^n$. Let $\vec{r} = (r_1, \dots, r_{Q/2})$.
 - Generates the punctured key $DK' \leftarrow \text{Puncture}(DK, \{c_1, r_1, \dots, c_{Q/2}, r_{Q/2}\})$
 - Samples $b \leftarrow \{0, 1\}$ and sends the following to \mathcal{A} :

$$\begin{array}{ll} (\vec{c}, \vec{r}, EK, DK') & \text{if } b = 0 \\ (\vec{r}, \vec{c}, EK, DK') & \text{if } b = 1 \end{array}$$

3. The adversary outputs b' and wins if $b = b'$.

¹⁵This Q will correspond to the number of queries made by the `Verify` algorithm in the final watermarking construction, hence the choice of letter.

¹⁶We can assume that the set C is represented as a list in sorted order.

Sparseness We also require that most strings are not valid ciphertexts:

$$\Pr \left[\text{Dec}(DK, c) \neq \perp \mid (EK, DK) \leftarrow \text{Gen}(1^\lambda), c \leftarrow \{0, 1\}^n \right] \leq \text{negl}(\lambda)$$

One of our contributions is the following theorem.

Theorem A.2. *A puncturable encryption system can be constructed using indistinguishability obfuscation and injective one-way functions*

A full construction and proof is provided in appendices B and B.2.

B Puncturable Encryption Construction

We provide a construction of the puncturable encryption defined in Appendix A.

B.1 Construction

We construct a puncturable encryption scheme in which the length n of ciphertexts is 12 times the length ℓ of plaintexts. Our construction utilizes the following ingredients:

- A length-doubling $PRG : \{0, 1\}^\ell \rightarrow \{0, 1\}^{2\ell}$
- A puncturable, injective family of PRFs $\{\mathcal{F}_\lambda : \{0, 1\}^{3\ell} \rightarrow \{0, 1\}^{9\ell}\}$. We require \mathcal{F}_λ to be selectively puncturable on any Q prefixes.¹⁸ [TODO: Remark that this can be done from any OWF.]
- A puncturable family of PRFs $\{\mathcal{G}_\lambda : \{0, 1\}^{9\ell} \rightarrow \{0, 1\}^\ell\}$. We require \mathcal{G}_λ to be selectively puncturable on any Q of points. [TODO: Remark that this can be done from any OWF.]
- A injective bit commitment Commit using randomness in $\{0, 1\}^{9\ell}$, which can in fact be constructed by an injective one-way function. We only use this in our security proof.

Construction B.1 (Puncturable Encryption Scheme). [TODO: Padding]

- $\text{Gen}(1^\lambda)$ samples a function $F \leftarrow \mathcal{F}$ and $G \leftarrow \mathcal{G}$, and generates EK as the iO-obfuscation of the circuit in Figure 4. Gen returns (EK, DK) , where DK is the (un-obfuscated) program in Figure 5.
- $\text{Puncture}(VK, c_0, c_1)$ outputs DK' , where DK' is the iO-obfuscation of the program described in Figure 6.
- $\text{Enc}(EK, m)$ takes $m \in \{0, 1\}^\ell$ and outputs $EK(m) \in \{0, 1\}^{12\ell}$.
- $\text{Dec}(DK, c)$ takes $c \in \{0, 1\}^{12\ell}$ and returns $DK(c)$.

Remark 4. We note that in all of our obfuscated programs (including the hybrids), whenever α_0 and α_1 or β_0 and β_1 or γ_0 and γ_1 are treated symmetrically, then we can and do store them in lexicographical order. A random ordering would also suffice for security.

¹⁸As in [SW14], any puncturable PRF family from $\{0, 1\}^k \rightarrow \{0, 1\}^{2k+\omega(\log \lambda)}$ can be made statistically injective (with no additional assumptions) by utilizing a family of pairwise-independent hash functions.

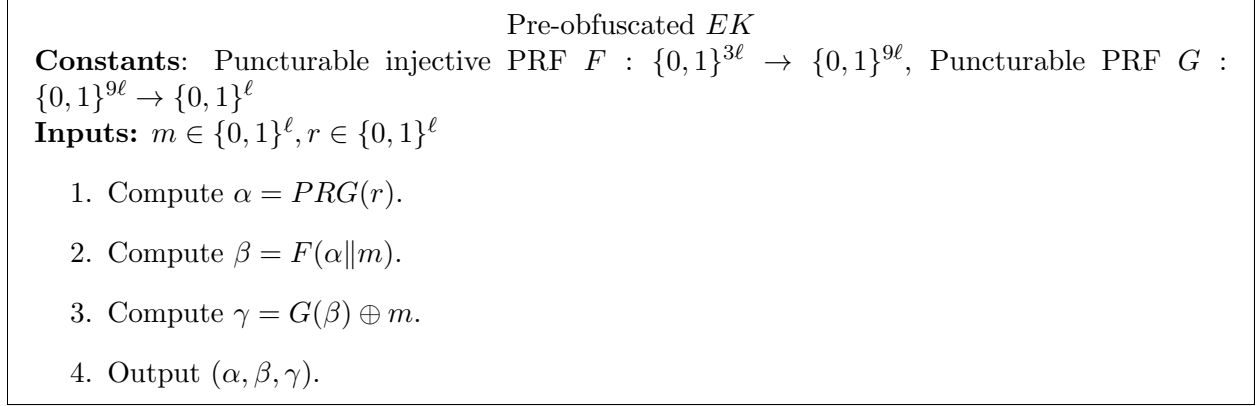


Figure 4: Program describing how to encode

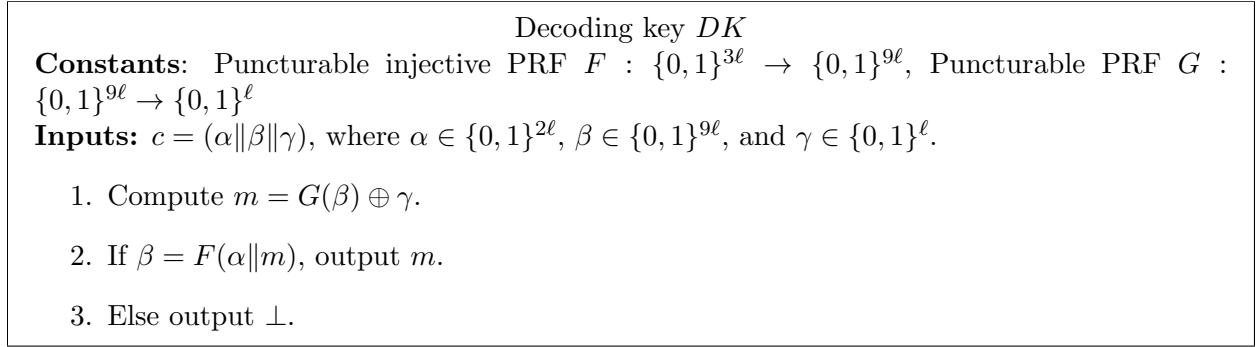


Figure 5: Program describing how to decode

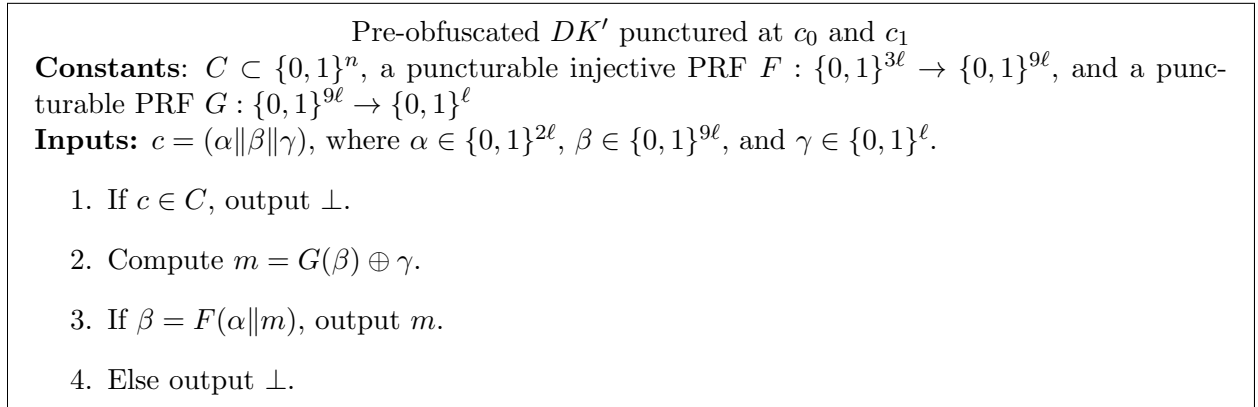


Figure 6: Program describing how to decode

Correctness Correctness follows from the fact that indistinguishability obfuscation exactly preserves functionality, and observing in the punctured case that DK' is defined to be functionally

equivalent to DK except on inputs in C .

Sparseness Sparseness follows from, for example, the length-doubling PRG; most values of α are not in the image of PRG .

B.2 Ciphertext Pseudorandomness

We give a sequence of hybrids H_0 through H_{14} . The goal of the hybrids to reach a game in which the challenge encryptions $c_1, \dots, c_{Q/2}$ and the random ciphertexts $r_1, \dots, r_{Q/2}$ are treated symmetrically in EK and DK' , and in which both are sampled uniformly at random by the challenger. We proceed by iteratively replacing pieces of c_0 by uniformly random values, puncturing F and G as necessary.

[TODO: Instead of simply presenting the hybrids in order, do a “top-down” view of the hybrid argument. Present the first/last hybrids, then the main 2-3 intermediate hybrids, then the rest of the intermediate hybrids.]

H_0 Hybrid H_0 is defined as the real security game:

1. \mathcal{A} sends a messages $m_1, \dots, m_{Q/2} \in \mathcal{M}$ to the challenger.
2. The challenger does the following:
 - (a) Samples an injective PRF $F : \{0, 1\}^{3\ell} \rightarrow \{0, 1\}^{9\ell}$ which is selectively puncturable on Q prefixes, and PRF $G : \{0, 1\}^{9\ell} \rightarrow \{0, 1\}^\ell$ selectively puncturable on Q points.
 - (b) For each $i \in [Q/2]$:
 - Samples $t_i \leftarrow \{0, 1\}^\ell$,
 - $\alpha_i = PRG(t_i) \in \{0, 1\}^{2\ell}$,
 - $\beta_i = F(\alpha_i \| m_i)$,
 - $\gamma_i = G(\beta_i) \oplus m_i$.
 - Let c_i as $\alpha_i \| \beta_i \| \gamma_i$, and $\vec{c} = (c_1, \dots, c_{Q/2})$.
 - (c) Samples $r_1, \dots, r_{Q/2} \leftarrow \{0, 1\}^{12\ell}$.
 - Parse $r_i = \alpha'_i \| \beta'_i \| \gamma'_i$ and let $\vec{r} = (r_1, \dots, r_{Q/2})$.
 - (d) Generates EK as the iO-obfuscation of [Figure 4](#) and DK' as the iO-obfuscation of [Figure 6](#).
 - (e) Samples $b \leftarrow \{0, 1\}$ and sends the following to \mathcal{A} :

$$\begin{aligned} (\vec{c}, \vec{r}, EK, DK') & \text{ if } b = 0 \\ (\vec{r}, \vec{c}, EK, DK') & \text{ if } b = 1 \end{aligned}$$

3. The adversary outputs b' and wins if $b = b'$.

H_1 : In hybrid H_1 , we alter the generation of the challenge ciphertexts (see Line 2(b) of H_0). Sample each $\alpha_i \leftarrow \{0, 1\}^{2\ell}$ uniformly at random.

H_2 : In hybrid H_2 , we alter the generation of EK (see Line 2(d) of H_0). We puncture F in EK on all strings of the form $\alpha_i \| \star$ or $\alpha'_i \| \star$ for each $i \in [Q/2]$. This is functionally equivalent because F is never evaluated on strings of these forms because α_i and α'_i are with high probability not in the image of PRG . This is where we use the the prefix-puncturability of F .

H_3 : In hybrid H_3 , we modify the generation of DK' . For each $i \in [Q/2]$, hard-code the constants $\hat{\beta}_i = F(\alpha'_i \| m_i)$ and $\hat{\gamma}_i = G(\hat{\beta}_i) \oplus m_i$. For each $i \in [Q/2]$, add the following line in the beginning of DK' : “If $c \in \alpha'_i \| \hat{\beta}_i \| \hat{\gamma}_i$, output m_i .” This change is functionally equivalent, as $\alpha_i \| \hat{\beta}_i \| \hat{\gamma}_i$ is already a valid encryption of c_i . Notice, that these $\hat{\beta}_i$ do not correspond to either the β_i or β'_i (and similarly for $\hat{\gamma}_i$). **[TODO: explain what they are. Forward pointer to H4.]**

For reference, we describe DK' from Hybrid H_3 in [Figure 7](#).

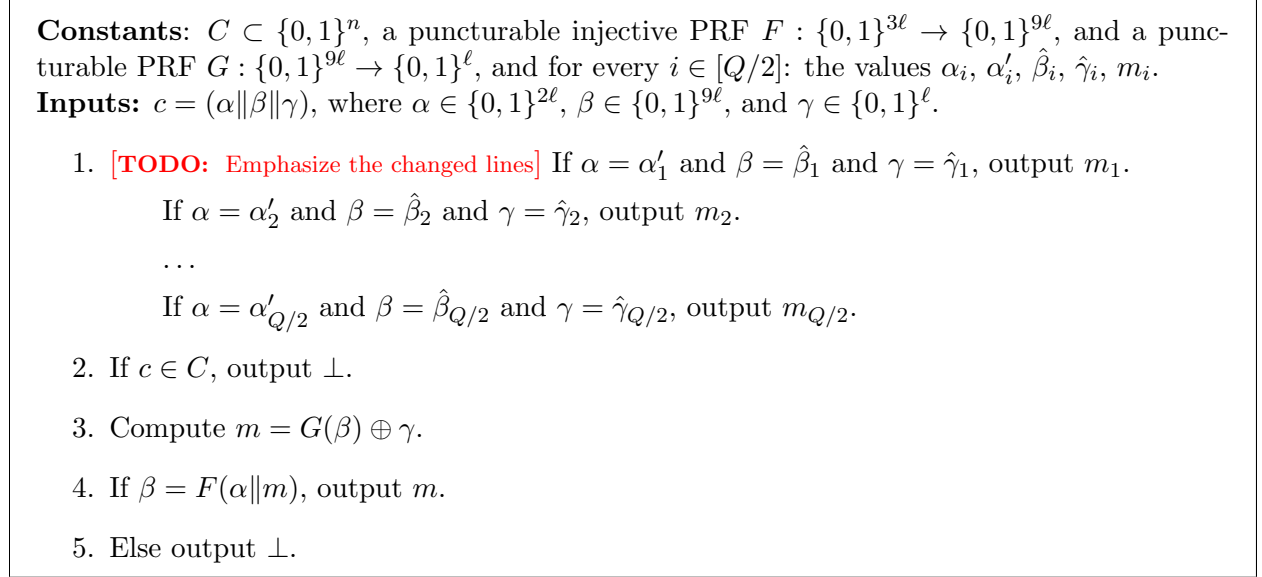


Figure 7: DK' as in Hybrid H_3 , pre-obfuscation

H_4 : In hybrid H_4 , we again modify the generation of DK' (see [Figure 8](#)). We add the following check: “If $(\alpha, m) \in \bigcup_{i=1}^{Q/2} \{(\alpha_i, m_i), (\alpha'_i, m_i)\}$, output \perp .” This is functionally equivalent by two cases:

1. When $(\alpha, m) = (\alpha_i, m_i)$ for some $i \in [Q/2]$, then either $c = c_i$, in which case DK' already would output \perp , or $c \neq c_i$, in which case DK' rejects c as an invalid ciphertext (because every pair (α, m) together define a unique valid ciphertext).
2. When $(\alpha, m) = (\alpha'_i, m_i)$, we only reach this line if $c \neq \alpha'_i \| \hat{\beta}_i \| \hat{\gamma}_i$ (by the check introduced in Hybrid H_3). In this case, DK' already rejects c as an invalid ciphertext.

For reference, we describe DK' from Hybrid H_4 in [Figure 8](#).

Constants: $C \subset \{0, 1\}^n$, a puncturable injective PRF $F : \{0, 1\}^{3\ell} \rightarrow \{0, 1\}^{9\ell}$, and a puncturable PRF $G : \{0, 1\}^{9\ell} \rightarrow \{0, 1\}^\ell$, and for every $i \in [Q/2]$: the values $\alpha_i, \alpha'_i, \hat{\beta}_i, \hat{\gamma}_i, m_i$.
Inputs: $c = (\alpha \parallel \beta \parallel \gamma)$, where $\alpha \in \{0, 1\}^{2\ell}$, $\beta \in \{0, 1\}^{9\ell}$, and $\gamma \in \{0, 1\}^\ell$.

1. If $\alpha = \alpha'_1$ and $\beta = \hat{\beta}_1$ and $\gamma = \hat{\gamma}_1$, output m_1 .
 If $\alpha = \alpha'_2$ and $\beta = \hat{\beta}_2$ and $\gamma = \hat{\gamma}_2$, output m_2 .
 ...
 If $\alpha = \alpha'_{Q/2}$ and $\beta = \hat{\beta}_{Q/2}$ and $\gamma = \hat{\gamma}_{Q/2}$, output $m_{Q/2}$.
2. If $c \in C$, output \perp .
3. Compute $m = G(\beta) \oplus \gamma$.
4. **[TODO: emphasize the changed lines.]** If $(\alpha, m) \in \bigcup_{i=1}^{Q/2} \{(\alpha_i, m_i), (\alpha'_i, m_i)\}$, output \perp .
5. If $\beta = F(\alpha \parallel m)$, output m .
6. Else output \perp .

Figure 8: DK' as in Hybrid H_4 , pre-obfuscation

H_5 : In hybrid H_5 , we modify the generation of the key DK' in the security game. Instead of using the unpunctured key for F , we puncture F at the points $\alpha_i \parallel m_i$ and $\alpha'_i \parallel m_i$ for each $i \in [Q/2]$. This is functionally equivalent because – by the checks added in the previous hybrid – F will never be evaluated on such inputs.

H_6 : In hybrid H_6 , we alter the generation of the challenge ciphertexts (see Line 2(b) of H_0). We sample β^* uniformly at random from $\{0, 1\}^{9\ell}$. This is indistinguishable by the pseudorandomness of F at punctured points.

H_7 : In hybrid H_7 , we change Line 1 of Figure 8. For each $i \in [Q/2]$, hardcode $z_i := \text{Commit}(0; \hat{\beta}_i)$, and we replace the check “ $\beta = \hat{\beta}_i$ ” with the check “ $\text{Commit}(0; \beta) = z_i$ ”. This is functionally equivalent by the injectivity of Commit .

H_8 : In hybrid H_8, z_i is instead hard-coded as “ $\text{Commit}(1; \hat{\beta}_i)$ ”. Indistinguishability is by the computational hiding of Commit .

H_9 : In hybrid H_9 , we replace the expression “ $\text{Commit}(0; \beta) = z_i$ ” with FALSE. This is functionally equivalent with high probability because of the perfect binding of Commit (which follows from injectivity). In fact, we remove the entire line 1, which also preserves functionality.

For reference, we describe DK' from Hybrid H_9 in Figure 10.

[TODO: Punctured keys] [TODO: update] Constants: $C \subset \{0, 1\}^n$, a punctured injective PRF $F' : \{0, 1\}^{3\ell} \rightarrow \{0, 1\}^{9\ell}$ punctured at $\bigcup_{i=1}^{Q/2} \{\alpha_i \| m_i, \alpha'_i \| m_i\}$, and a puncturable PRF $G' : \{0, 1\}^{9\ell} \rightarrow \{0, 1\}^\ell$. Also the values α_i, α'_i, m_i for each $i \in [Q/2]$. **Inputs:** $c = (\alpha \| \beta \| \gamma)$, where $\alpha \in \{0, 1\}^{2\ell}$, $\beta \in \{0, 1\}^{9\ell}$, and $\gamma \in \{0, 1\}^\ell$.

1. If $(\alpha, m) \in \bigcup_{i=1}^{Q/2} \{(\alpha_i, m_i), (\alpha'_i, m_i)\}$, output \perp .
2. If $c \in C$, output \perp .
3. Compute $m = G(\beta) \oplus \gamma$.
4. If $\beta = F'(\alpha \| m)$, output m .
5. Else output \perp .

Figure 9: DK' as in Hybrid H_4 , pre-obfuscation

H_{10} : In hybrid H_{10} , we modify how the challenge ciphertexts are generated (see Line 2(b) of H_0). For every $i \in [Q/2]$, sample $\beta_i \leftarrow \{0, 1\}^{9\ell}$ uniformly at random. This is indistinguishable by the pseudorandomness of F at the (selectively) punctured points.

H_{11} : In hybrid H_{11} , we alter the generation of EK (see Line 2(d) of H_0). We puncture G in EK on β_i and β'_i for every $i \in [Q/2]$. This is functionally equivalent by the sparsity of F ; since β_i and β'_i are now chosen at random for every i , with high probability they are not in the image of F .

H_{12} : In hybrid H_{12} , we alter the generation of DK' , changing Line 2 of Figure 10. Instead of “If $c \in C$: output \perp ”, we replace it with “If $\beta \in \bigcup_{i=1}^{Q/2} \{\beta_i, \beta'_i\}$: output \perp ”. To see that this change is functionally equivalent, we observe that with high probability, neither of these lines has any effect.

Since with high probability, none of the β_i and β'_i are in the image of F , if $\beta \in \bigcup_{i=1}^{Q/2} \{\beta_i, \beta'_i\}$ – which is the case when $c \in C$ – then $DK'(c) = \perp$ with high probability, even without the extra check.

The obvious question, then, is why not remove the check? Because checking if $\beta \in \bigcup_{i=1}^{Q/2} \{\beta_i, \beta'_i\}$ will allow us to puncture G on this set in the following hybrid.

H_{13} : In hybrid H_{13} , we alter the generation of DK' . We puncture G at $\bigcup_{i=1}^{Q/2} \{\beta_i, \beta'_i\}$ in DK' . This change is functionally equivalent because of the ostensibly useless checks in the previous hybrid.

H_{14} : In hybrid H_{14} , we generate the challenge ciphertexts in a different way. For each $i \in [Q/2]$, we sample γ_i uniformly at random from $\{0, 1\}^\ell$. This change is indistinguishable by the selective indistinguishability of G at the punctured set.

For reference, we describe DK' from Hybrid H_{14} in Figure 10. In this hybrid, $c_i = \alpha_i \| \beta_i \| \gamma_i$ and $r_i = \alpha'_i \| \beta'_i \| \gamma'_i$ are now treated symmetrically in both EK and DK' . Furthermore, they are both sampled uniformly and independently at random from $\{0, 1\}^{12\ell}$. So no adversary has any advantage in this hybrid.

Constants: $C \subset \{0, 1\}^n$, a punctured injective PRF $F' : \{0, 1\}^{3\ell} \rightarrow \{0, 1\}^{9\ell}$ punctured at $\bigcup_{i=1}^{Q/2} \{\alpha_i \| m_i, \alpha'_i \| m_i\}$, and a punctured PRF $G' : \{0, 1\}^{9\ell} \rightarrow \{0, 1\}^\ell$ punctured at $\bigcup_{i=1}^{Q/2} \{\beta_i, \beta'_i\}$. Also the values α_i, α'_i, m_i for each $i \in [Q/2]$.

Inputs: $c = (\alpha \| \beta \| \gamma)$, where $\alpha \in \{0, 1\}^{2\ell}$, $\beta \in \{0, 1\}^{9\ell}$, and $\gamma \in \{0, 1\}^\ell$.

1. If $(\alpha, m) \in \bigcup_{i=1}^{Q/2} \{(\alpha_i, m_i), (\alpha'_i, m_i)\}$, output \perp .
2. If $c \in C$, output \perp .
3. Compute $m = G'(\beta) \oplus \gamma$.
4. If $\beta = F'(\alpha \| m)$, output m .
5. Else output \perp .

Figure 10: DK' as in Hybrid H_{14} , pre-obfuscation

C Proof of Theorem 5.3

To prove Theorem 5.3, we first give a distinguishing game and show that no p.p.t. algorithm \mathcal{A} can win this game with non-negligible probability. Because Verify is black-box as in Definition 4.1 and queries exactly 2 points of an argument program, a variant of Theorem 4.2 implies the desired result.¹⁹

Game C.1 (Lunchtime Distinguishing). First, the challenger generates (mk, vk) by $\text{Setup}(1^\lambda)$. The adversary is presented with vk and access to the following two oracles.

- A marking oracle \mathcal{O}_M which takes a circuit C and returns $\text{Mark}(\text{mk}, C)$.
- A challenge oracle \mathcal{O}_C which takes no input, but samples a circuit C^* uniformly from \mathcal{C} and returns $C^*_\# = \text{Mark}(\text{mk}, C^*)$. The adversary must query \mathcal{O}_C exactly once, after which he may no longer query either oracle.

At the end of the game, the challenger executes $\text{Verify}(\text{vk}, C^*_\#)$. From the points at which Verify queried $C^*_\#$, the challenger picks x_0 uniformly at random. The challenger also chooses x_1 as a random point in the domain of $C^*_\#$. The challenger picks a random bit b and sends x_b to the adversary.

The adversary then outputs a bit b' and wins if $b' = b$. The adversary's advantage is $|\Pr[b' = b] - \frac{1}{2}|$.

Lemma C.2. *Every p.p.t. algorithm \mathcal{A} has negligible advantage in Theorem C.1 for Theorem 5.2.*

We will show that security in the last hybrid reduces to the security of the puncturable encryption \mathcal{PE} . We will construct an algorithm \mathcal{B} distinguishing $(EK, DK\{x_0, x_1\}, x_0, x_1)$ from $(EK, DK\{x_0, x_1\}, x_1, x_0)$ with non-negligible advantage. To do so, we will have to answer the queries of the watermarking adversary in Theorem C.1 using the \mathcal{PE} challenge. The main challenge therefore is to use only the punctured decryption key and to treat x_0 and x_1 symmetrically.

¹⁹As noted in the discussion of Theorem 4.2, we prove that lemma for specific removing and distinguishing games. The proof techniques and result apply directly in this setting as well.

Proof. We must show that in the game, no p.p.t. algorithm can distinguish x_0 from a random point x_1 with non-negligible advantage. Recall that x_0 is chosen from the points at which `Verify` queries $C_{\#}^*$. In [Theorem 5.2](#), there are two such points, one of which is $PRG_2(PRG_1(a^*))$ for randomly chosen $a^* \leftarrow \{0, 1\}^{\ell/4}$, and thus pseudorandom. In the rest of the proof, we focus on the other case: x_0 is computed as $\text{Enc}(EK, PRG_1(a^*) || h^*)$, where $h^* = h(C_{\#}^*(PRG_2(PRG_1(a^*))))$. We must show that this distribution of x_0 is indistinguishable from a uniform x_1 .

As \mathcal{A} 's runtime is bounded by some polynomial q , we define a sequence of hybrid games H_0 through H_{q+4} . Hybrid H_0 is simply [Theorem C.1](#), except that the challenger picks C^* , generates $C_{\#}^*$, and samples x_0 and x_1 appropriately before presenting the adversary with vk . This is possible because these are all chosen independently of the adversary's actions.

In hybrid H_i , the adversary's first i queries C_1, \dots, C_i to \mathcal{O}_m are answered in a modified way (as in [Figure 11](#)). These queries are answered using only a *punctured* decryption key $DK' = DK\{x_0, x_1\}$ and *punctured* PRF key $F' = F\{x_0, x_1\}$. In hybrid H_q , we answer all queries to \mathcal{O}_M in this way.²⁰

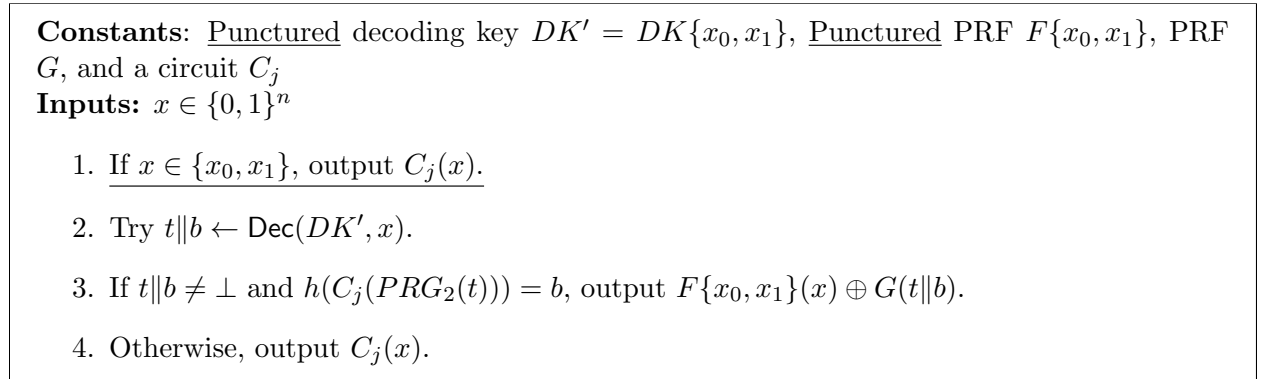


Figure 11: Modified marked program $C_{j\#}$ for $j \leq i$ in hybrid H_i .

Claim C.2.1. *In H_i , $h(C_{i+1}(PRG_2(PRG_1(a^*)))) = h^*$ with negligible probability.*

Proof. First we note that because $h : \{0, 1\}^m \rightarrow \{0, 1\}^{\ell/2}$ is a compressing collision-resistant hash function, h is also a one-way function. Suppose that in hybrid H_i , $h(C_{i+1}(PRG_2(PRG_1(a^*)))) = h^*$ with non-negligible probability. Then we provide an algorithm `Inv` that inverts h , or violates the pseudorandomness of C^* . Throughout this proof, let PRG denote $PRG_2 \circ PRG_1$.

`Inv` takes as input $y = h(r)$ for uniformly sampled $r \leftarrow \{0, 1\}^m$. `Inv` samples random $a^* \leftarrow \{0, 1\}^{\ell/2}$, $EK, DK \leftarrow \text{Gen}(1^\lambda)$, $x_1 \leftarrow \{0, 1\}^n$, $F \leftarrow \mathcal{F}_\lambda$, and $G \leftarrow \mathcal{G}_\lambda$. `Inv` samples an encryption of $a^* || y$, namely $x_0 \leftarrow \text{Enc}(a^* || y)$. Lastly, `Inv` punctures both the decryption key and PRF key at $\{x_0, x_1\}$, obtaining DK' and F' .

Using the above, `Inv` runs \mathcal{A} and answers queries to \mathcal{O}_M as in hybrid H_i . The view of \mathcal{A} in this simulation is indistinguishable from the view in the real hybrid H_i ; the only difference is that in the real hybrid, $y = h(r)$ for $r = C^*(PRG(a^*))$ for random $C^* \leftarrow \mathcal{C}$ (instead of uniformly random r). By the pseudorandomness of the family \mathcal{C} , the views are indistinguishable.

²⁰The reason for puncturing F here is so that later we may make a change to the challenge program (see hybrid H_{q+3}).

Finally, \mathcal{A} outputs a circuit C_{i+1} . By our hypothesis, $h(C_{i+1}(PRG(a^*))) = h^*$ with non-negligible probability. Thus, with non-negligible probability, Inv has found a pre-image of y under h , namely $C_{i+1}(PRG(a^*))$ \square \square

The proof of the above claim crucially relies on the fact that when \mathcal{A} queries the circuit C_i , he has no information about C^* . This is why this proof only applies to *lunchtime* unremovability. Using the above claim, we now prove the following:

Claim C.2.2. For $0 \leq i < q$, $H_i \approx H_{i+1}$.

Proof. The only difference between H_i and H_{i+1} is in the way that the $i + 1^{\text{th}}$ query is answered. We now show that the two programs returned are functionally equivalent with high probability, and so by the security of iO , the two hybrids are indistinguishable.

By the correctness of the punctured \mathcal{PE} decryption key DK' and the correctness of the punctured PRF F' on non-punctured points, there are only two possible inputs on which $C_{i+1\#}$ may differ in H_i and H_{i+1} : namely, x_0 and x_1 .

By the sparseness of ciphertexts in \mathcal{PE} , $\text{Dec}(DK, x_1) = \perp$ with high probability over the choice of x_1 . Thus in hybrid H_i , $C_{i+1\#}(x_1) = C_{i+1}(x_1)$. This is also true in hybrid H_{i+1} .

On the other hand, x_0 decrypts to $a^* || h^*$. By the previous claim, the check that $h(C_{i+1}(PRG(a^*))) = h^*$ fails with high probability. Thus in hybrid H_i , $C_{i+1\#}(x_0) = C_{i+1}(x_0)$. This is also true in hybrid H_{i+1} . \square \square

From H_q to H_{q+4} , every query to \mathcal{O}_M is answered as [Figure 11](#). We proceed to modify only $C_{\#}^*$, with two simultaneous goals in mind. We want x_0 and x_1 to be treated symmetrically in the *challenge* marked program $C_{\#}^*$. We also want to generate $C_{\#}^*$ using only the *punctured* decryption key $DK'\{x_0, x_1\}$.

H_{q+1} : In hybrid H_{q+1} , we modify $C_{\#}^*$ as in [Figure 12](#). We puncture C^* at $\{x_1\}$ and F at $\{x_0\}$, and hard-code a mapping $x_0 \mapsto y_0$ and $x_1 \mapsto y_1$. y_0 and y_1 are defined as $y_0 = F(x_0) \oplus G(PR_{G_1}(a^*) || h^*)$ and $y_1 = C^*(x_1)$. We also puncture the decryption key DK at $\{x_0, x_1\}$. These changes preserve functionality with high probability. F is never evaluated on x_0 because of the hard-coded check in line 1. Similarly, $C^*\{x_1\}$ is not evaluated at x_1 on line 4. Furthermore, with high probability, x_1 is not in the image of PRG . Thus on line 3, $C^*\{x_1\}$ is never evaluated at x_1 .

Constants: A punctured decoding key $DK' = DK\{x_0, x_1\}$, punctured PRF $F\{x_0\}$, PRF G , and a punctured circuit $C^*\{x_1\}$, values $x_0, x_1, y_0 = F(x_0) \oplus G(PR_{G_1}(a^*) || h^*), y_1 = C^*(x_1)$

Inputs: $x \in \{0, 1\}^n$

1. If $x = x_i$ for $i \in \{0, 1\}$, output y_i .
2. Try $a || b \leftarrow \text{Dec}(DK', x)$.
3. If $a || b \neq \perp$ and $h(C^*\{x_1\}(PRG_2(PR_{G_1}(a)))) = b$, output $F\{x_0\}(x) \oplus G(PR_{G_1}(a) || b)$.
4. Otherwise, output $C^*\{x_1\}(x)$.

Figure 12: Modified $C_{\#}^*$ in H_{q+1}

H_{q+2} : In hybrid H_{q+2} , we change y_1 to be random. This is an indistinguishable change by the punctured pseudorandomness of C^* at x_1 .

H_{q+3} : In hybrid H_{q+3} , we change y_0 to be random. This is an indistinguishable change by the punctured pseudorandomness of F at x_0 .

H_{q+4} : In hybrid H_{q+4} , we no longer puncture C^* and F . This preserves the functionality of $C_{\#}^*$ by the same reasoning as in hybrid H_{q+1} , hence is an indistinguishable change by the security of iO .

Final Security Argument We now reduce to the security of \mathcal{PE} . Given an adversary \mathcal{A} with non-negligible advantage in hybrid H_{q+4} , we construct an adversary \mathcal{B} violating the punctured ciphertext indistinguishability of \mathcal{PE} . That is, in [Theorem A.1](#), \mathcal{B} distinguishes $(EK, DK\{x_0, x_1\}, x_0, x_1)$ from $(EK, DK\{x_0, x_1\}, x_1, x_0)$ with the same advantage. \mathcal{B} executes the following steps:

1. \mathcal{B} samples $C^* \leftarrow \mathcal{C}$, picks $a^* \leftarrow \{0, 1\}^{\ell/4}$, computes $h^* = h(C^*(PRG_2(PRG_1(a^*))))$, and sends $m = a^* \| h^*$ to the \mathcal{PE} challenger.
2. \mathcal{B} receives $(EK, DK\{x_0, x_1\}, x_b, x_{1-b})$ as input, where $x_0 \leftarrow \text{Enc}(m)$, $x_1 \leftarrow \{0, 1\}^n$ for some unknown $b \in \{0, 1\}$.
3. \mathcal{B} samples $F \leftarrow \mathcal{F}$ and $G \leftarrow \mathcal{G}$ to construct the verification vk as in [Setup](#), and runs \mathcal{A} on vk . \mathcal{B} answers \mathcal{A} 's queries to \mathcal{O}_M and \mathcal{O}_C as in H_{q+4} .
4. At the end of the game, \mathcal{B} sends x_b to \mathcal{A} . \mathcal{A} outputs a bit b' , and \mathcal{B} also outputs b' .

In this execution, \mathcal{A} 's view is exactly the same as in H_{q+4} . The b in [Theorem C.1](#) is the same as the b in [Theorem A.1](#). \square

D Proof of [Theorem 5.4](#)

For clarity, we recall the security game of (q, γ) -relaxed unforgeability. In the theorem statement, $q = 1 - \gamma$.

Game D.1 ((q, γ) -relaxed Unforgeability). First, the challenger generates $(\text{mk}, \text{vk}) \leftarrow \text{Setup}(1^\lambda)$. The adversary is presented with vk and access to a marking oracle \mathcal{O}_M , which takes a circuit C_i and returns $C_{i\#} = \text{Mark}(\text{mk}, C_i)$.

At the end of the game, the adversary outputs a circuit \hat{C} . The adversary wins if the following conditions hold:

1. $\Pr \left[\text{Verify}(\text{vk}, \hat{C}) = 1 \right] \geq q + \frac{1}{\text{poly}(\lambda)}$
2. $\Pr \left[\forall i, \hat{C}(x) \neq C_i(x) \mid x \leftarrow \{0, 1\}^n \right] \geq \gamma$.

Remark 5. As in the proof of [Theorem 4.2](#), we will analyze \mathcal{A} 's winning probability conditioned on \mathcal{A} outputting a \hat{C} that satisfies condition 2. This isn't an efficiently testable event because the

threshold of γ is too sharp.²¹ However, as discussed in [Theorem 4.2](#), we can *relax* this condition in a way which only increases \mathcal{A} 's advantage. Even in this game, the adversary has a negligible chance of winning. For clarity, we omit these technical details and just assume that \mathcal{A} always outputs such a \hat{C} .

of [Theorem 5.4](#). We will construct a series of games H_0 through H_3 between a challenger and an adversary. Each game defines some random variables. In particular, each game defines the following variables:

- The queries C_1, \dots, C_T made by the adversary (T is a bound on the adversary's running time)
- The candidate forgery \hat{C} output by the adversary at the end of the game
- A randomly chosen t^* . This, together with C^*, C_1, \dots, C_T defines h_1, \dots, h_T and h^* . Specifically, we define $h_i = h(C_i(\text{PRG}_2(t^*)))$, and $h^* = h(\hat{C}(\text{PRG}_2(t^*)))$.
- x^*

In our security proof, we consider two events:

1. The event that $\hat{C}(x^*) = F(x^*) \oplus G(t^* \| h^*)$. In H_0 , the probability of this event is exactly the probability that $\text{Verify}(\text{vk}, \hat{C}) = 1$ in [Theorem D.1](#).
2. The event that for some $i \in \{1, \dots, T\}$, $h_i = h^*$. We will upper bound the probability of this “bad” event to upper bound the probability that $\text{Verify}(\text{vk}, \hat{C}) = 1$.

H_0 : Our first hybrid H_0 , is defined as follows: the challenger runs [Theorem D.1](#) with [Theorem 5.2](#), as well as sampling a few extra variables. The extra variables are $a^* \leftarrow \{0, 1\}^{\ell/4}$ and $t^* = \text{PRG}_1(a^*)$.

That is, the challenger first samples $(EK, DK) \leftarrow \mathcal{PE.Gen}(1^\lambda)$, and samples $F \leftarrow \mathcal{F}$ and $G \leftarrow \mathcal{G}$. These are used to define marking keys (mk, vk) as in $\text{Setup}(1^\lambda)$. The challenger sends vk to \mathcal{A} .

When \mathcal{A} makes a query C_i , the challenger responds with $C_{i\#} = \text{Mark}(\text{mk}, C_i)$. If the adversary produces a candidate forgery \hat{C} , x^* is defined as $\text{Enc}(EK, t^* \| h^*)$, where $h^* = h(\hat{C}(\text{PRG}_2(t^*)))$.

Claim D.1.1. *In H_0 , event 2 happens with probability at most $1 - \gamma + \text{negl}(\lambda)$.*

Proof. Here we just use the collision resistance of h , the pseudorandomness of PRG_2 , and the fact that in [Theorem D.1](#), the adversary is restricted to producing \hat{C} such that \hat{C} differs from *every* C_i on at least a γ fraction of the domain. \square

Claim D.1.2. *In H_0 , event 1 happens with probability at most $1 - \gamma + \text{negl}(\lambda)$*

Proof. We give indistinguishable hybrids H_1 through H_3 , such that in hybrid H_3 , the probability of event 1 is bounded by the probability of event 2, which by the previous claim is $1 - \gamma + \text{negl}(\lambda)$.

H_1 : In hybrid H_1 , \mathcal{B} samples $t^* \leftarrow \{0, 1\}^{\ell/2}$ instead of generating $t^* = \text{PRG}_1(a^*)$. This is indistinguishable by the pseudorandomness of PRG_1 .

H_2 : In hybrid H_2 , \mathcal{B} sends a modified vk to the adversary. In this vk , the puncturable PRF G is punctured at the set of strings beginning with t^* . This is functionally equivalent because with high

²¹in fact, $\#P$ -complete!

probability t^* is not in the image of PRG_1 . Indistinguishability thus follows from the security of iO .

H_3 : In hybrid H_3 , \mathcal{B} answers the queries to \mathcal{O}_M in a modified way. In particular, each marked response $C_{i\#}$ will have G punctured on the set of strings beginning with t^* . $C_{i\#}$ will also have the values t^* , h_i , and $G(t^*||h_i)$ hard-coded where $h_i = h(C_i(PRG_2(t^*)))$. We modify $C_{i\#}$ by the correct hard-coded value on the only punctured point ($t^*||h_i$) on which G will ever be evaluated. Indistinguishability thus follows from the security of iO .

When event 2 does not happen, the probability that $\hat{C}(x^*) = F(x^*) \oplus G(t^*||h^*)$ is negligible. \square

This concludes the proof of [Theorem 5.4](#). \square \square

E Proof of [Theorem 4.2](#)

of [Theorem 4.2](#). Suppose that a removing adversary outputs a circuit \hat{C} such that $\Pr[\text{Verify}(\text{vk}, \hat{C}) = 0] = L \cdot \delta + \text{Adv}$ for some non-negligible Adv . Let q_i be a random variable denoting the i^{th} point at which Verify queries $C_{\#}^*$.

$$\begin{aligned} L \cdot \delta + \text{Adv} &= \Pr[\text{Verify}(\text{vk}, \hat{C}) = 0] \\ &\leq \Pr[\exists i : \hat{C}(q_i) \neq C_{\#}^*(q_i)] + \text{negl}(\lambda) \end{aligned} \tag{1}$$

$$\begin{aligned} &\leq \sum_i \Pr[\hat{C}(q_i) \neq C_{\#}^*(q_i)] + \text{negl}(\lambda) \tag{2} \\ &= L \cdot \Pr_{i \leftarrow \{1, \dots, L\}} [\hat{C}(q_i) \neq C_{\#}^*(q_i)] + \text{negl}(\lambda) \\ &= L \cdot \Pr_{x_0} [\hat{C}(x_0) \neq C_{\#}^*(x_0)] + \text{negl}(\lambda) \end{aligned}$$

Inequality (1) is by the black-box property of Verify . With high probability, $\text{Verify}(\text{vk}, C_{\#}^*) = 1$, by completeness. If $\text{Verify}(\text{vk}, \hat{C}) = 0$, then (with high probability) there must be some queried point q_i for which $\hat{C}(q_i) \neq C_{\#}^*(q_i)$. Inequality (2) is by a union bound. As a result $\Pr_{x_0} [\hat{C}(x_0) \neq C_{\#}^*(x_0)] \geq \delta + \text{Adv}'$ for some non-negligible $\text{Adv}' \approx \frac{\text{Adv}}{L}$.

If the adversary outputs a circuit \hat{C} such that $\hat{C} \sim_{\delta} C_{\#}^*$, then since x_1 is chosen uniformly at random, $\Pr_{x_1} [\hat{C}(x_1) \neq C_{\#}^*(x_1)] \leq \delta$. Thus there is a p.p.t. algorithm distinguishing x_0 from x_1 with non-negligible advantage (specifically, $\text{Adv}'/2$). \square \square

F Unwatermarkable PRFs

Our starting point is the constructions of unobfuscatable function families in [\[BGI⁺12\]](#) and [\[BP12\]](#), and an understanding of those constructions will prove helpful towards understanding ours. The former work presents a construction of 0-robustly extractable PRF families; from any exact implementation, the key can be recovered. They extend this notion to a very weak form of approximate functionality (weaker than what we require). The latter work handles a very strong form of approximation: the approximate implementation must only agree on some constant fraction of the

domain. They achieve this, they sacrifice the total learnability of the earlier construction, instead learning only a single predicate of the PRF key. We require a notion of approximation stronger than [BGI⁺12] but weaker than [BP12], and a notion of learnability weaker than [BGI⁺12] but stronger than [BP12], and achieve this by adapting techniques from both works.

F.1 Preliminaries

The construction requires an invoker randomizable pseudorandom function [BGI⁺12] and a decomposable encryption schemes [BP12]. The following definitions and discussion are taken almost verbatim from those works.

Definition F.1 (Invoker-Randomizable Pseudorandom Functions, [BGI⁺12]). A function ensemble $\{f_k\}_{k \in \{0,1\}^*}$ such that $f_k : \{0,1\}^{n+m} \rightarrow \{0,1\}^m$, where n and m are polynomially related to $|k|$, is called an *invoker-randomizable pseudorandom function ensemble* if the following holds:

1. $\{f_k\}_{k \in \{0,1\}^*}$ is a PRF family.
2. For every k and $x \in \{0,1\}^n$, the mapping $r \mapsto f_k(x, r)$ is a permutation over $\{0,1\}^m$.

Property 2 implies that, for every fixed k and $x \in \{0,1\}^n$, if r is chosen uniformly in $\{0,1\}^m$, then the value $f_k(x, r)$ is distributed uniformly (and independently of x) in $\{0,1\}^m$.

Lemma F.1 ([BGI⁺12]). *If pseudorandom functions exist, then there exist invoker-randomizable pseudorandom functions.*

Definition F.2 (Decomposable Encryption [BP12]). An encryption scheme $(\text{Gen}, \text{Enc}, \text{Dec})$ is *decomposable* if there exists an efficient algorithm pub that operates on ciphertexts and satisfies the following conditions:

1. For a ciphertext c , $\text{pub}(c)$ is independent of the plaintext and samplable; that is, there exists an efficient sampler PubSamp such that, for any secret key $sk \in \{0,1\}^n$:

$$\text{PubSamp}(1^n) \equiv \text{pub}(\text{Enc}_{sk}(0)) \equiv \text{pub}(\text{Enc}_{sk}(1))$$

2. A ciphertext c is deterministically defined by $\text{pub}(c)$ and the plaintext; that is, for every secret key sk and two distinct ciphertexts c and c' , if $\text{pub}(c) = \text{pub}(c')$, then $\text{Dec}_{sk}(c) \neq \text{Dec}_{sk}(c')$.

We use as our decomposable encryption scheme a specific symmetric-key encryption scheme which enjoys a number of other necessary properties. Given a PRF $\{f_k\}_{k \in \{0,1\}^*}$ with one-bit output and for security parameter λ , the secret key is a random $sk \in \{0,1\}^\lambda$, and the encryption of a bit b is computed by sampling a random $r \leftarrow \{0,1\}^\lambda$ and outputting $(r, F_{sk}(r) \oplus b)$. This function satisfies a number of necessary properties [BP12]:

- It is CCA-1 secure.
- It is decomposable.
- The support of $(\text{Enc}_{sk}(0))$ and $(\text{Enc}_{sk}(1))$ are each a non-negligible fraction (in reality, at least $\frac{1}{2} - \text{negl}$) of the cipher-text space.
- For a fixed secret key sk , random samples from $(b, \text{Enc}_{sk}(b))_{b \leftarrow \{0,1\}}$ are indistinguishable from uniformly random strings.

F.2 Construction

The key k for the PRF is given by a tuple $k = (\alpha, \beta, sk, s_1, s_2, s_e, s_h, s_b, s^*)$. For security parameter λ , α and β are uniformly random λ -bit strings, sk is a secret key for the decomposable encryption scheme described above, s_h is a key for an invoker-randomizable pseudorandom function, and s_1, s_2, s_e, s_b , and s^* are independent keys for a family of PRFs. We denote by F_s a PRF with key s .

The domain of the PRF will be of the form (i, q) for $i \in \{1, \dots, 9\}$, and $q \in \{0, 1\}^{\ell(n)}$, for some polynomial ℓ . The range is similarly bit strings of length polynomial in ℓ . The function will be defined in terms of 9 auxiliary functions, and the index i will select among them. We use a combination of ideas from [BGI⁺12] and [BP12] to construct a PRF family for which s^* can be recovered from any (negligibly-close) approximation to f_k , which will enable us to compute f_k restricted to $i = 9$. This allows us to recover a $1/9$ -close approximation of f_k that is implementation independent (simply by returning 0 whenever $i \neq 9$). To achieve a δ -close approximation for any $\delta = \frac{1}{\text{poly}(\lambda)}$, we simply augment the index i with an additional $\log(\delta)$ bits: if all these bits are 0, then we index as before; otherwise, use index $i = 9$. Instead of recovering $1/9$ th of the function, we now recover $1 - \delta$ of the function. This establishes the theorem.²²

We now define the auxiliary functionalities we will use in the construction.

- \mathbb{R}_s : The function \mathbb{R}_s is parameterized by a PRF key s . It takes as input q and returns $\mathbb{R}_s(q) = F_s(q)$, the PRF evaluated at q . That is, \mathbb{R}_s simply evaluates a PRF.
- $\mathbb{C}_{a,b,s}$: The function $\mathbb{C}_{a,b,s}$ is parameterized by two bit strings a and b , and a PRF key s . It takes as input q and returns $\mathbb{C}_{a,b,s}(q) = b \oplus F_s(q \oplus a)$, where F_s is the PRF given by key s . That is, \mathbb{C} evaluates a PRF on a point related to the queried point, then uses the value to mask the bitstring b .
- $\mathbb{E}_{sk,\alpha,s_e}$: The function $\mathbb{E}_{sk,\alpha,s_e}$ is parameterized by a secret key sk for the encryption scheme, a bitstring α , and a PRF key s_e . It takes as input q and returns $\mathbb{E}_{sk,\alpha,s_e}(q) = \text{Enc}_{sk}(\alpha; r)$ with randomness $r = F_{s_e}(q)$. That is, \mathbb{E} returns an encryption of α using randomness derived by evaluating the PRF on the query.
- \mathbb{H}_{sk,s_h} : The function \mathbb{H}_{sk,s_h} is parameterized by a secret key sk for the encryption scheme, and an invoker-randomizable PRF key s_h . It takes as input two cipher-texts of bits c and d , the description of a two-bit gate \odot , and some additional input \bar{q} , and returns $\mathbb{H}_{sk,s_h}(c, d, \odot, \bar{q}) = \text{Enc}_{sk}(\text{Dec}_{sk}(c) \odot \text{Dec}_{sk}(d); r)$ with randomness $r = F_{s_h}(c, d, \odot, \bar{q})$. That is, \mathbb{H} implements a homomorphic evaluation of \odot on the ciphertexts c and d by decrypting and re-encrypting, with randomness derived by applying a PRF to the whole input.
- $\mathbb{B}_{sk,\alpha,\beta,s_b}$: The function $\mathbb{B}_{sk,\alpha,\beta,s_b}$ is parameterized by a secret key sk for the symmetric-key encryption scheme, bitstrings α and β , and a PRF key s_b . It takes as input n ciphertexts c_1, \dots, c_λ and additional input \bar{q} , and returns

$$\mathbb{B}_{sk,\alpha,\beta,s_b}(c_1, \dots, c_\lambda, \bar{q}) = \alpha \oplus F_{s_b}(m_1 \oplus \beta_1, \dots, m_\lambda \oplus \beta_\lambda, \text{pub}(c_1), \dots, \text{pub}(c_\lambda), \bar{q})$$

²²Note that the result is a PRF family that depends on the choice of δ . The argument would fail if δ was a negligible function, because an approximation for could “erase” all the structure of the PRF family, thwarting learnability. Removing this dependence (ie: constructing a family that works for all inverse polynomial δ simultaneously) would be interesting.

where $m_i = \text{Dec}_{sk}(c_i)$.

Having defined the auxiliary functions, our pseudorandom function f_k for $k = (\alpha, \beta, sk, s_1, s_2, s_e, s_h, s_b, s^*)$ is a combination of these functions. The argument (i, q) selects which function is evaluated, and q is parsed appropriately by each of the functionalities. For example, \mathbb{B} parses q as λ ciphertexts c_1, \dots, c_λ , and all remaining bits as \bar{q} .

$$f_k(i, q) = \begin{cases} \mathbb{C}_1(q) := \mathbb{C}_{\alpha, \beta, s_1}(q) & \text{if } i = 1 \\ \mathbb{C}_2(q) := \mathbb{C}_{\alpha, s^*, s_2}(q) & \text{if } i = 2 \\ \mathbb{E}(q) := \mathbb{E}_{sk, \alpha, s_e}(q) & \text{if } i = 3 \\ \mathbb{H}(q) := \mathbb{H}_{sk, s_h}(q) & \text{if } i = 4 \\ \mathbb{B}(q) := \mathbb{B}_{sk, \alpha, \beta, s_b}(q) & \text{if } i = 5 \\ \mathbb{R}_1 := \mathbb{R}_{s_1}(q) & \text{if } i = 6 \\ \mathbb{R}_2 := \mathbb{R}_{s_2}(q) & \text{if } i = 7 \\ \mathbb{R}_b := \mathbb{R}_{s_b}(q) & \text{if } i = 8 \\ \mathbb{R}^* := \mathbb{R}_{s^*}(q) & \text{if } i = 9 \end{cases}$$

While this construction may appear daunting, each subfunction serves a very concrete purpose in the argument; understanding the proof ideas will help clarify the construction. We must now argue two properties of this family: learnability as in [Definition 6.2](#), and pseudorandomness.

F.3 Learnability

We must show that $F_\lambda = \{f_k\}$ is robustly, $\frac{1}{9}$ -approximately learnable by an implementation-independent algorithm, L .²³ It suffices to show that, given any ρ -implementation g of f_k for random key k , s^* can be recovered, because $\mathbb{R}^* = \mathbb{R}_{s^*}$ comprises 1/9th of the functionality.

To begin, consider the case when the implementation is perfect: $g \equiv f_k$. In this case, recovery of s^* is straightforward. Given α , \mathbb{C}_1 , and \mathbb{R}_1 it is easy to find β : for any q , $\beta = \mathbb{C}_1(q) \oplus \mathbb{R}_1(q \oplus \alpha)$. That is, it is easy to construct a circuit that, on input α , outputs β (by fixing some uniformly random q in the above).²⁴ But we don't know α , only encryptions of α (coming from \mathbb{E}), so how might we recover β ?

Using \mathbb{H} , it is easy to homomorphically evaluate the circuit on such an encryption, yielding an encryption $c = (c_1, \dots, c_n)$ of $\beta = (\beta_1, \dots, \beta_n)$. For any \bar{q} , evaluating $\mathbb{B}(c, \bar{q})$ will yield $\alpha \oplus F_{s_b}(\vec{0}, c, \bar{q})$. Evaluating $\mathbb{R}_b(\vec{0}, \text{pub}(c_1), \dots, \text{pub}(c_n), \bar{q})$ immediately yields α in the clear. Now we can directly recover $s^* = \mathbb{C}(q) \oplus \mathbb{R}_2(q \oplus \alpha)$, for any q .

How does this argument change when g and f_k may disagree on an (arbitrary) ρ -fraction of the domain for some negligible function $\rho(n)$? The first observation is that in the above algorithm, each of \mathbb{C}_1 , \mathbb{C}_2 , \mathbb{E} , \mathbb{R}_1 , and \mathbb{R}_2 , can each be evaluated (homomorphically in the case of \mathbb{C}_1) at a single point that is distributed uniformly at random. With high probability, g will agree with f_k on these inputs.

²³As discussed earlier, it suffices to prove learnability for $\delta = 1/9$. We may then change how the subfunctions are indexed to achieve any inverse polynomial.

²⁴This ability is what enables the learnability; the black-box learner cannot construct such a circuit and thus cannot continue with the homomorphic evaluation in the next step.

It remains to consider robustness to error in \mathbb{H} , \mathbb{B} , and \mathbb{R}_b . The same idea does not immediately work, because the queries to these circuits are not uniform.

For \mathbb{H} , we leverage the invoker-randomizability of the PRF F_{s_h} , using the argument presented in [BGI⁺12]²⁵. In every query to $\mathbb{H}(c, d, \odot, \bar{q})$, the input \bar{q} only effects the randomness used in the final encrypted output. For each such query, pick \bar{q} uniformly and independently at random. Now \mathbb{H} returns a uniformly random encryption of $\text{Dec}_{sk}(c) \odot \text{Dec}_{sk}(d)$. This is because the randomness used for the encryption is now uniformly sampled by F_{s_h} . The distribution over the output induced by the random choice of \bar{q} depends only on $(\text{Dec}_{sk}(c), \text{Dec}_{sk}(d), \odot) \in \{0, 1\}^2 \times \{0, 1\}^2 \times \{0, 1\}^4$. As in [BGI⁺12], the probability of returning an incorrect answer on such a query is at most 64ρ , which is still negligible.

For \mathbb{B} and \mathbb{R}_b , we leverage the properties of the decomposable symmetric-key encryption scheme, using the argument presented in [BP12].²⁶ We modify the procedure of using \mathbb{B} and \mathbb{R}_b to recover α given an encryption c of β . Instead of querying \mathbb{B} on (c, \bar{q}) , sample a fresh random m , and using \mathbb{H} , compute an encryption c' of $\beta \oplus m$. Note that c' is a uniformly random encryption (by invoker-pseudorandomness) of the uniformly random string $\beta \oplus m$, and is thus a uniformly-distributed string of the appropriate length. Independently sample a random \bar{q} and query $\alpha' := \mathbb{B}(c', \bar{q})$. This query to \mathbb{B} is now distributed uniformly, and will therefore be answered correctly with high probability.

To recover α , we evaluate $\alpha = \alpha' \oplus \mathbb{R}_b(m, \text{pub}(c_1), \dots, \text{pub}(c_\lambda), \bar{q})$. This query to \mathbb{R}_b is also distributed uniformly at random (for random \bar{q}), and will therefore be answered correctly with high probability.

F.4 Pseudorandomness

Our proof that the family $\{f_k\}$ is pseudorandom follows that of [BP12]; the main technical change comes from the fact that \mathbb{B} depends on α . We consider a polynomial-time adversary \mathcal{A} with oracle access to f_k . For simplicity, we ignore the indexing of the subfunctions of f_k and assume that \mathcal{A} has direct oracle access to each of the constituent functions, showing that they are simultaneously pseudorandom.

Let E_1 be the the event that \mathcal{A} produces *distinct* queries $q = (c, \bar{q})$, $q' = (c', \bar{q}')$ such that:

$$(m \oplus \beta, \text{pub}(c_1), \dots, \text{pub}(c_\lambda), \bar{q}) = (m' \oplus \beta, \text{pub}(c'_1), \dots, \text{pub}(c'_\lambda), \bar{q}')$$

where $m, m' \in \{0, 1\}^\lambda$ are the decryptions under sk of c and c' respectively.

Claim F.1.1. $\Pr_{k, \mathcal{A}}[E_1] = 0$

Proof. Recall that for any ciphertext c , $\text{pub}(c)$ and the plaintext m uniquely determine the ciphertext. If $m \oplus \beta = m' \oplus \beta$, and $\text{pub}(c_i) = \text{pub}(c'_i)$ for all i , then $c = c'$. Therefore $q = q'$. \square

We consider two “bad” events, and argue that if \mathcal{A} is to distinguish f_k from a random function, (at least) one of the events must occur.

- Let E_α be the event that \mathcal{A} produces queries q and q' such that $q \oplus \alpha = q'$.
- Let E_β be the event that \mathcal{A} produces queries $q = (c, \bar{q})$ and q' such that $q' = (m \oplus \beta, \text{pub}(c_1), \dots, \text{pub}(c_\lambda), \bar{q})$, where $m \in \{0, 1\}^\lambda$ is the decryption under sk of c .

²⁵Proof of Theorem 4.3

²⁶Proof of Claim 3.8

Claim F.1.2. *If $\Pr_{k,\mathcal{A}}[E_\alpha] \leq \text{negl}(\lambda)$ and $\Pr_{k,\mathcal{A}}[E_\beta] \leq \text{negl}(n)$, then \mathcal{A} cannot distinguish between f_k and a random oracle.*

Proof. Because f_k depends on the PRF keys $s_1, s_2, s_e, s_h,$ and s_b (but not s^*) only by black-box application of the respective PRFs, we can indistinguishably replace all applications of these PRFs by (independent) truly random functions. If E_α never occurs, then the responses from \mathbb{C}_1 and \mathbb{R}_1 (respectively \mathbb{C}_2 and \mathbb{R}_2) are uncorrelated; thus we can indistinguishably replace \mathbb{C}_1 (respectively, \mathbb{C}_2) by a independent random function. At this point, \mathcal{A} 's oracle only depends on s^* through calls to the PRF F_s^* ; we can now replace \mathbb{R}^* with a independent random function. By similar reasoning, if E_β never occurs, then the responses from \mathbb{B} and \mathbb{R}_b are uncorrelated; thus we can indistinguishably replace \mathbb{B} with another independent random oracle. The above holds with high probability, conditioning on $\neg E_\alpha$ and $\neg E_\beta$.

Now \mathcal{A} is left with oracles of \mathbb{E} and \mathbb{H} in which the PRFs F_{s_e} and F_{s_h} have been replaced by random (along with 7 additional independent random oracles). The ciphertexts of the encryption scheme we use are pseudorandom. Thus, access to these two oracles may be replaced with random without noticeably affecting the output distribution of \mathcal{A} . \square

All that remains is to bound the probabilities of E_α and E_β . We consider two cases separately: when E_α occurs before E_β and vice-versa, arguing that the probability of either event occurring first is negligible. Let $E_{\alpha,i}$ (respectively, $E_{\beta,i}$) be the event that E_α (respectively E_β) occurs in the first i queries.

Claim F.1.3. *For all i , $\Pr_{k,\mathcal{A}}[E_{\beta,i} | \neg E_{\alpha,i-1}] \leq \text{negl}(\lambda)$*

Proof. It suffices to show that for all i :

$$\Pr_{k,\mathcal{A}}[E_{\beta,i} | \neg E_{\alpha,i-1}, \neg E_{\beta,i-1}] \leq \text{negl}(\lambda).$$

Furthermore, because the events are efficiently testable given only $\alpha, \beta,$ and sk , it is enough to prove the claim when all the underlying PRFs (corresponding to $s_1, s_2, s_e, s_h, s_b,$ and s^* are replaced by (independent) truly random functions.

As in Claim F.1.2, if E_α doesn't occur in the first $i - 1$ queries, then the responses from \mathbb{C}_1 and \mathbb{R}_1 (respectively \mathbb{C}_2 and \mathbb{R}_2) are uncorrelated on these queries; thus we can indistinguishably replace \mathbb{C}_1 (respectively, \mathbb{C}_2) by a independent random function. By similar reasoning, if E_β doesn't occur in the first $i - 1$ queries, then the responses from \mathbb{B} and \mathbb{R}_b are uncorrelated on these queries; thus we can indistinguishably replace \mathbb{B} with another independent random oracle. The above holds with high probability, conditioning on $\neg E_{\alpha,i-1}$ and $\neg E_{\beta,i-1}$.

The view of \mathcal{A} after the first $i - 1$ queries is now independent of β . Now E_β amounts to outputting a ciphertext c and string q such that $\text{Dec}_{sk}(c) \oplus q = \beta$, for $\beta \leftarrow \{0, 1\}^\lambda$ drawn independently of the view of the adversary. This occurs with vanishingly small probability. \square

Claim F.1.4. $\Pr_{k,\mathcal{A}}[E_{\alpha,i} | \neg E_{\beta,i-1}] \leq \text{negl}(\lambda)$

Proof. It suffices to show that for all i :

$$\Pr_{k,\mathcal{A}}[E_{\alpha,i} | \neg E_{\beta,i-1}, \neg E_{\alpha,i-1}] \leq \text{negl}(\lambda).$$

Again, because the events are efficiently testable given only $\alpha, \beta,$ and sk , it is enough to prove the claim when all the underlying PRFs (corresponding to $s_1, s_2, s_e, s_h, s_b,$ and s^* are replaced by

(independent) truly random functions. As in the previous claim, we may indistinguishably replace the first i - responses of $\mathbb{C}_1, \mathbb{C}_2, \mathbb{B}, \mathbb{R}_b, \mathbb{R}_1,$ and \mathbb{R}_2 by independent random functions. The above holds with high probability, conditioning on $\neg E_{\alpha, i-1}$ and $\neg E_{\beta, i-1}$.

The view of the adversary is depends on α only by way of \mathbb{E} , the circuit that outputs random encryptions of α . Furthermore, besides the oracles \mathbb{E} and \mathbb{H} , all of the oracle responses \mathcal{A} receives are uniformly random (and independent of α). But just as in [BGI⁺12]²⁷ and [BP12]²⁸, with only these two oracles, any CCA-1 encryption scheme is semantically secure. Thus we can indistinguishably replace $\mathbb{E}_{sk, \alpha, s_e}$ with $\mathbb{E}_{sk, \alpha, s_e} -$ returning only encryptions of 0. Finally, the view of \mathcal{A} is information theoretically independent of α ; as before, we conclude that $E_{\alpha, i}$ occurs with vanishingly small probability. \square

G Relationship between γ and δ

Below we illustrate some simple requirements on δ and γ in the unremovability and unforgeability definitions that are necessary if both are to be satisfied simultaneously. Specifically, if a watermarking scheme is both $(0, \delta)$ -unremovable and $(0, \gamma)$ -unforgeable, then $\gamma \geq \delta + \frac{1}{\text{poly}(n)}$ for some polynomial poly.

Recall that an adversary δ -removes for a challenge $C_{\#}^*$ if it outputs a program $\hat{C} \approx_{\delta} C_{\#}^*$ such that $\text{Verify}(VK, \hat{C}) = 0$ with non-negligible probability. That is, it must remove the mark without changing the challenge program on more than δ -fraction of inputs. An adversary wins the γ -forges if it outputs a program \hat{C} such that $\text{Verify}(VK, \hat{C}) = 1$ with non-negligible probability and additionally, for all marked programs $C_{i\#}$ seen by the adversary, $\hat{C} \not\approx_{\gamma} C_{i\#}$. That is, its output is only considered a forgery if it is at least γ -far from all marked programs.

Consider the following ‘‘attack’’: Given a random marked program $C_{\#} : \{0, 1\}^n \rightarrow \{0, 1\}^m$, consider the following program, parameterized by $c \in [2^n]$:

$$C_c(x) = \begin{cases} C_{\#}(x) \oplus 1 & \text{if } x \leq c \\ C_{\#}(x) & \text{if } x > c \end{cases}$$

Consider $b \leftarrow \text{Verify}(C_c, VK)$. At least one of $\Pr[b = 0]$ or $\Pr[b = 1]$ is at least $1/2$. In the former case, this construction violates unremovability unless $\delta \leq \frac{c}{2^n}$; in the latter case, this construction violates unforgeability unless $\gamma \geq \frac{c}{2^n}$. That is, if a watermarking scheme is both unremovable and unforgeable for parameters δ and γ , then for all $c \in [2^n]$, either $\gamma \geq \frac{c}{2^n}$ or $\delta \leq \frac{c}{2^n}$. Therefore, $\gamma \geq \delta$.

Furthermore, in the setting when Verify is black-box with respect to C_c , then $\gamma \geq \delta + \frac{1}{\text{poly}(n)}$, for some polynomial. We alter the above attack as follows. For $c \in [2^n]$ and a pseudo-random permutation $\pi \leftarrow \mathcal{PRP}$:

$$P_{c, \pi}(x) = \begin{cases} C_{\#}(x) \oplus 1 & \text{if } \pi(x) \leq c \\ C_{\#}(x) & \text{if } \pi(x) > c \end{cases}$$

Suppose there is some negligible function $\text{negl}(n)$ such that $\gamma = \delta + \text{negl}(n)$. For $c = \delta 2^n$ and randomly chosen $\pi \leftarrow \mathcal{PRP}$, $\text{Verify}(VK, C_{\delta 2^n, \pi}) = 1$ with all but negligible probability; otherwise $C_{\delta 2^n, \pi}$ violates unremovability. For $c = \gamma 2^n$, $\text{Verify}(VK, C_{\gamma 2^n, \pi}) = 0$ with all but negligible probability; otherwise $C_{\gamma 2^n, \pi}$ violates unforgeability.

²⁷Claim 3.6.1

²⁸Claim 3.3

The programs $C_{\delta 2^n, \pi}$ and $C_{\gamma 2^n, \pi}$ disagree on a negligible fraction of the domain. The set on which they disagree is pseudo-random, by the security of π . Because `Verify` is black-box, we can use it to distinguish black-box access to these two functions, which is statistically-difficult.

G.1 Multi-bit Equivalence

In our main theorem, we state that there exist watermarking schemes for puncturable pseudorandom function families with long output lengths (specifically, $\Omega(\lambda^\epsilon)$ for some constant $\epsilon > 0$).

However, we note that there is a simple reduction, allowing us to can watermark puncturable pseudorandom function families with arbitrary output lengths, suffering only a small loss in parameters.

The idea is simple: any pseudorandom function with multi-bit outputs can be viewed as a PRF with single-bit outputs, with a slightly expanded input space is slightly expanded.

Concretely, if we have a (p, δ) -unremovable watermarking scheme for a PRF family \mathcal{F} with m -bit outputs, then this is easily seen to be a $(p, \delta/m)$ -unremovable watermarking scheme for the \mathcal{F} interpreted as a PRF family with single-bit outputs. If the scheme for \mathcal{F} is (q, γ) -unforgeable, then it is also (q, γ) -unforgeable for \mathcal{F} when interpreted as a family with single-bit outputs.

Dually, given a scheme for a family with single-bit outputs, we can analyze its parameters when construed as a scheme for an m -bit PRF family. A (p, δ) -unremovable scheme for a single-bit family is also (p, δ) -unremovable for the same family construed as an m -bit family. A (q, γ) -unforgeable scheme for the single-bit family becomes a $(q, m\gamma)$ -unforgeable scheme.