

# Program Obfuscation II: The Puncturing Paradigm

## 1 IO is a “Best-Possible” Obfuscation

Assume that in the future someone discovers an obfuscator  $\mathcal{O}'$  which has better hiding properties, e.g. given  $\mathcal{O}'(C)$ , it is hard for any PPT adversary to predict  $\pi(C)$  for a class of predicates  $\pi$ . Assume that  $\mathcal{O}'$  has overhead  $q(\cdot)$  for some polynomial  $q$  in the sense that  $|\mathcal{O}'(C)| \leq q(|C|)$ .

We will show that  $\mathcal{O}(\text{pad}_q(C))$ , where  $\mathcal{O}$  is an indistinguishability obfuscator, has the *same* hiding properties where  $\text{pad}_q(C)$  makes the circuit  $C$  artificially bigger, of size  $q(|C|)$ . Assume, for contradiction, that there is an adversary  $\mathcal{A}$  such that

$$\Pr[\mathcal{A}(\mathcal{O}(\text{pad}_q(C))) = \pi(C)] \geq 1/2 + 1/p(\lambda).$$

Since  $\text{pad}_q(C)$  and  $\mathcal{O}'(C)$  compute the same function and are of the same size,

$$\Pr[\mathcal{A}(\mathcal{O}(\mathcal{O}'(C))) = \pi(C)] \geq 1/2 + 1/p(\lambda) - \text{negl}(\lambda)$$

by the indistinguishability property of  $\mathcal{O}$ . This gives us an adversary who contradicts the claimed security of  $\mathcal{O}'$ . The adversary  $\mathcal{A}'$ , on input  $\mathcal{O}'(C)$ , simply runs  $\mathcal{A}$  on  $\mathcal{O}(\mathcal{O}'(C))$ , giving us

$$\Pr[\mathcal{A}'(\mathcal{O}'(C)) = \pi(C)] \geq 1/2 + 1/p(\lambda) - \text{negl}(\lambda),$$

contradicting the claimed security of  $\mathcal{O}'$ .

### 1.1 An Application: Time-lock Puzzles

A time-lock puzzle is defined by a puzzle generation algorithm  $\text{Gen}$  which outputs a puzzle  $k$  together with its solution  $s$ . We have an  $\text{Eval}$  algorithm that is given  $k$ , runs in *sequential* time  $T$  and outputs  $s$ . We require that no PPT machine that runs in *sequential time* significantly less than  $T$  should be able to output  $s$  given  $k$ , even if it can do arbitrary polynomial-time parallel computation  $\gg T$ .

A classical example is the Rivest-Shamir-Wagner construction where the puzzle is  $(g, x, N)$ :  $N = pq$  is an RSA modulus;  $g$  is a random element of  $\mathbb{Z}_N^*$ , and  $x \in_R [0, N]$  is randomly chosen. The solution is  $s = g^{x^T} \bmod N$ . This can be computed in time proportional to  $T$  (think of  $T$  as being much larger than  $\text{poly}(\log N)$ , the time to do basic operations in the group). The point is,  $s$  is easy to compute given the factorization of  $N$  which in particular tells us  $\varphi(N)$ . Given this, first compute  $y = x^T \bmod \varphi(N)$  and then raise  $g$  to the power of  $y \bmod N$ . However, given only  $N$ , and assuming factoring is hard, it appears that any strategy would require computing iteratively  $g^{x^i} \bmod N$  for  $i = 1, 2, \dots, T$ .

Are there different constructions? In particular, constructions that don't have to make the seemingly ad-hoc assumption that exponentiation is inherently sequential?

**A Construction using IO.** If you think about it a bit, you realize that one has to assume *something*. i.e. that there exist computations that are inherently sequential. More formally, there is a Turing machine  $M_T$  such that on input  $x$ , computing  $M_T(x)$  requires *sequential time*  $T$ . In other words, we have to assume that not all computations are parallelizable.

Consider a dummy Turing machine  $M'_{T,s}$  that does nothing for  $T$  time steps, and eventually outputs the string  $s$ . The description of the puzzle is the IO obfuscation  $k = \mathcal{O}(M'_{T,s})$ . Here, we assume that one can in fact obfuscate Turing machines in a way that the size of the obfuscated TM does not depend on the runtime

of the TM. This can indeed be done assuming IO for circuits; see Canetti-Jain-Holmgren-Vaikuntanathan STOC 2015.

We claim that no  $T' \ll T$  time algorithm can learn  $s$  given the puzzle  $k$ . Suppose, for contradiction, that there is an algorithm that does this. Consider running the algorithm on input  $\mathcal{O}(M''_{T,s})$  where  $M''_{T,s}$  runs the inherently sequential machine  $M_T$  for  $T$  time steps, outputs  $s$  if  $M_T$  outputs 1, and outputs  $\perp$  if  $M_T$  outputs 0.

On the one hand, if  $M_T$  outputs 1,  $M''_{T,s}$  and  $M'_{T,s}$  are functionally equivalent, so the adversary outputs  $s$  in this case, in time  $T' \ll T$ . On the other hand, if  $M_T$  outputs 0,  $M''_{T,s}$  is functionally equivalent to the machine that always outputs  $\perp$ . Here, the adversary has information-theoretically no chance of outputting  $s$  in polynomial-time. Thus, being able to output  $s$  is a distinguisher between the cases that  $M_T$  outputs 1 and where it outputs 0. In summary, an adversary that breaks the time-lock puzzle gives us an algorithm to compute  $M_T$  in sequential time  $T' \ll T$ , breaking the assumption that  $M_T$  takes sequential time  $T$  to compute.

## 2 IO Engineering I: Using Puncturable PRFs

Puncturable pseudorandom functions (PPRFs) are extremely useful in constructions using IO. We define them below and show several applications. Informally, PPRFs are PRFs where you can release keys that allows one to compute values of the PRF at all inputs except for one “punctured” point.

**Definition 1** (Puncturable PRFs). *A puncturable family of PRFs is given by a triple of uniform PPT machines (PRFGen, Punc, PRF), and a pair of computable functions  $n(\cdot)$  and  $m(\cdot)$ , satisfying the following conditions:*

**Correctness.** *For all outputs  $K$  of  $\text{PRFGen}(1^\lambda)$ , all points  $i \in \{0, 1\}^{n(\lambda)}$ , and  $K\{i\} = \text{Punc}(K, i)$ , we have that*

$$\text{PRF}(K\{i\}, x) = \text{PRF}(K, x)$$

*for all  $x \neq i$ .*

**Pseudorandom at punctured point.** *For every p.p.t. adversary  $(\mathcal{A}_1, \mathcal{A}_2)$ , there is a negligible function  $\mu$ , such that in an experiment where  $\mathcal{A}_1(1^\lambda)$  outputs a point  $i \in \{0, 1\}^{n(\lambda)}$  and a state  $\sigma$ ,  $K \leftarrow \text{PRFGen}(1^\lambda)$  and  $K\{i\} = \text{Punc}(K, i)$ , the following holds*

$$\left| \Pr[\mathcal{A}_2(\sigma, K\{i\}, i, \text{PRF}(K, i)) = 1] - \Pr[\mathcal{A}_2(\sigma, K\{i\}, i, U_{m(\lambda)}) = 1] \right| \leq \mu(\lambda)$$

*where  $\mu$  is called the distinguishing gap for  $(\mathcal{A}_1, \mathcal{A}_2)$ .*

*Furthermore, we say that the puncturable PRF is  $\delta$ -indistinguishable if the above pseudorandom property holds with a distinguishing gap  $\mu$  bounded by  $\delta$ . Especially, the puncturable PRF is sub-exponentially indistinguishable if  $\mu(\lambda)$  is bounded by  $2^{-\lambda^\epsilon}$  for a constant  $\epsilon$ .*

The GGM tree-based construction of PRFs from pseudorandom generators (PRGs) yields puncturable PRFs for any (efficiently computable)  $n(\cdot)$  and  $m(\cdot)$ . Furthermore, it is easy to see that if the PRG underlying the GGM construction is sub-exponentially hard (and this can in turn be built from sub-exponentially hard OWFs), then the resulting puncturable PRF is sub-exponentially pseudo-random.

## 2.1 Application: A Different Public-Key Encryption Scheme from IO

This is (likely?) the construction that Diffie and Hellman had in mind. The idea is to start from a secret-key encryption scheme, obfuscate its encryption algorithm, and publish the obfuscated program as the public key. More precisely, the program to be obfuscated contains the secret key  $sk$ , takes in a message  $m$  and a random string  $r$ , and outputs  $SKEnc(sk, m; r)$ . However, this has to be done with care.

As a first try, let us take the simple secret-key encryption scheme where the secret key  $sk$  is the seed of a pseudorandom function (PRF), and the encryption algorithm evaluates the PRF on a random string and uses it as a one-time pad for the message. That is,

$$SKEnc(sk, m; r) = (r, PRF(sk, r) \oplus m).$$

The public-key now is an obfuscated program that takes  $(m, r)$  as input and outputs  $SKEnc(sk, m; r)$ . Is the resulting public-key encryption scheme secure?

The construction that actually works is a clever twist on this, and is due to Sahai and Waters (Proceedings of ACM STOC 2014).

**The Construction.** Let  $\mathcal{O}$  be an indistinguishability obfuscator,  $G : \{0, 1\}^n \rightarrow \{0, 1\}^{2n}$  be a pseudorandom generator, and  $(PRFGen, Punc, PRF)$  be a puncturable PRF family. The public-key encryption scheme  $(KeyGen, Enc, Dec)$  works as follows.

- $KeyGen(1^n)$ : The secret key is a uniformly random string  $K \leftarrow \{0, 1\}^\lambda$ . Let  $\Pi_K$  be the following program that takes as input  $m \in \{0, 1\}^{2n}$  and  $r \in \{0, 1\}^n$ .

$\Pi_K(m, r)$  does the following:

- Compute  $y = G(r)$ .
- Output  $(y, PRF(K, y) \oplus m)$ .

The public key is an obfuscation of  $\Pi_K$  and the secret key is the PRF key  $K$ . That is,

$$pk = \hat{\Pi}_K = \mathcal{O}(\Pi_K) \quad \text{and} \quad sk = K.$$

- $Enc(pk, m)$ : Sample a random string  $r \leftarrow \{0, 1\}^n$  and output  $\hat{\Pi}_K(r, m)$ .
- $Dec(sk, c)$ : Run the secret-key decryption algorithm. That is, interpret  $c$  as a pair  $(c_0, c_1)$  and output  $c_1 \oplus PRF(K, c_0)$  as the message.

**Theorem 2.**  $(KeyGen, Enc, Dec)$  is a semantically secure public-key encryption scheme.

*Proof.* Correctness is immediate. We now show semantic security.

**Hybrid 0.** The real distribution of  $pk = \mathcal{O}(\Pi_K)$  and  $Enc(pk, m^*) = (y^* = G(r^*), PRF(K, y^*) \oplus m^*)$  for  $r^* \leftarrow \{0, 1\}^n$ .

**Hybrid 1.** We change the public key from  $\text{pk} = \mathcal{O}(\Pi_K)$  to  $\text{pk} = \mathcal{O}\left(\Pi'_{y^*, K\{y^*\}, \text{PRF}(K, y^*)}\right)$ , where  $\Pi'_{y^*, K\{y^*\}, z^*}(m, r)$  is the following program, taking as input  $m \in \{0, 1\}^{2n}$  and  $r \in \{0, 1\}^n$ .

$\Pi'_{y^*, K\{y^*\}, z^*}(m, r)$ :

- Compute  $y = G(r)$ .
- If  $y = y^*$ , output  $(y, z^* \oplus m)$ .
- Else, output  $(y, \text{PRF}(K\{y^*\}, y) \oplus m)$ .

Since we hard-code  $z^* = \text{PRF}(K, y^*)$  into the program,  $\Pi'_{y^*, K\{y^*\}, \text{PRF}(K, y^*)}$  has the exact same functionality as  $\Pi_K$ , and thus we can invoke IO security (as well as PPRF correctness) to argue indistinguishability from hybrid 0.

**Hybrid 2.** We substitute  $y^* = G(r^*)$  for  $y^* \leftarrow \{0, 1\}^{2n}$  in both the public key and ciphertext.

This is indistinguishable from hybrid 1 by security of the PRG. (The only reference to  $r^*$  is via  $y^* = G(r^*)$ , so this is valid.)

**Hybrid 3.** We change the public key into the obfuscation of the following program  $\Pi''_{y^*, K\{y^*\}}$ , taking as input  $m \in \{0, 1\}^{2n}$  and  $r \in \{0, 1\}^n$ .

$\Pi''_{y^*, K\{y^*\}}(m, r)$ :

- Compute  $y = G(r)$ .
- If  $y = y^*$ , output  $\perp$ .
- Else, output  $(y, \text{PRF}(K\{y^*\}, y) \oplus m)$ .

The only difference in functionality between  $\Pi''_{y^*, K\{y^*\}}$  and  $\Pi'_{y^*, K\{y^*\}, \text{PRF}(K, y^*)}$  is in the case when  $G(r) = y^*$ . However, since  $G$  is length-doubling and  $y^* \leftarrow \{0, 1\}^{2n}$  is uniformly random, the probability that there exists *any*  $r \in \{0, 1\}^n$  such that  $G(r) = y^*$  is at most  $2^{-n}$ . Therefore, with probability at least  $1 - 2^{-n}$  over  $y^*$ , the check " $G(r) \stackrel{?}{=} y^*$ " will never pass on any input  $r$ . Thus, with probability at least  $1 - 2^{-n}$ ,  $\Pi''_{y^*, K\{y^*\}}$  and  $\Pi'_{y^*, K\{y^*\}, \text{PRF}(K, y^*)}$  are functionally identical, allowing us to invoke IO security to argue indistinguishability from hybrid 2.

**Hybrid 4.** In the ciphertext, we replace  $\text{PRF}(K, y^*) \oplus m^*$  with a uniformly random string  $c^* \leftarrow \{0, 1\}^{2n}$ .

Since the public key now depends only on  $K\{y^*\}$  and not  $K$ , by PPRF security, the value  $\text{PRF}(K, y^*)$  is indistinguishable from random, and by a one-time pad argument, the same holds for  $c^*$ .

Hybrid 4 is independent of the message  $m^*$ , so we are done with the proof. □

### 3 When IO Alone is Not Enough

#### 3.1 Fully Homomorphic Encryption

A natural idea is to obfuscate a program that, on input any two ciphertexts  $C_0$  and  $C_1$  encrypting bits  $b_0$  and  $b_1$  respectively, decrypts them using a hardcoded secret-key  $SK$ , computes the NAND function on  $b_0$  and  $b_1$  and re-encrypts them under  $PK$ . That is,

$$\Pi_{PK,SK}(C_0, C_1) = \text{Enc}(PK, \text{Dec}(SK, C_0) \text{ NAND } \text{Dec}(SK, C_1); r)$$

An immediate question is where the randomness  $r$  for  $\text{Enc}$  comes from. It clearly shouldn't be hardcoded into  $\Pi$ ; neither should it be left up to the adversary to input the randomness into  $\Pi$ . The only option left is to hardcode a PRF key  $K$  into  $\Pi$ , and apply it to the input  $(C_0, C_1)$  to generate  $r$ . That is,

$$\Pi_{PK,SK,K}(C_0, C_1) = \text{Enc}(PK, \overline{\text{Dec}(SK, C_0) \wedge \text{Dec}(SK, C_1)}; \text{PRF}(K, (C_0, C_1)))$$

This is indeed the idea that will work out for us. In fact, it is not hard to see that if  $\mathcal{O}$  is a VBB obfuscation scheme, this construction gives us a secure FHE scheme. However, all we have is IO, and that requires quite a bit of care.

It turns out that any IO scheme plus a *perfectly rerandomizable encryption scheme* implies FHE. Both the construction and the proof is instructive and offers another technique to use IO in a world where IO and one-way functions together are not enough.

**The Construction.** Let's consider where the natural construction of obfuscating the circuit  $\Pi_{PK,SK,K}$  as above and releasing it publicly as the homomorphic evaluation key goes wrong. We'd like to argue that the evaluation key is indistinguishable from the obfuscation of a circuit

$$\Pi'_{PK,K}(C_0, C_1) = \text{Enc}(PK, 0; \text{PRF}(K, (C_0, C_1)))$$

which simply ignores the input (except to generate the encryption randomness) and outputs an encryption of 0. Note that this circuit does not use  $SK$  at all.

We will also make one more change: let the input and output encryption schemes use distinct keys. That is, we will have programs

$$\Pi_{PK_i,SK_{i-1},K_i}(C_0, C_1) = \text{Enc}(PK_i, \overline{\text{Dec}(SK_{i-1}, C_0) \wedge \text{Dec}(SK_{i-1}, C_1)}; \text{PRF}(K_i, (C_0, C_1)))$$

for  $i = 0, \dots, D - 1$  to evaluate circuits of depth up to  $D$ . We will show that this sequence of obfuscated programs is indistinguishable from the sequence of programs

$$\Pi'_{PK_i,K_i}(C_0, C_1) = \text{Enc}(PK_i, 0; \text{PRF}(K_i, (C_0, C_1)))$$

That is, we will show that

$$\left( \mathcal{O}(\Pi_{PK_1,SK_0,K_1}), \dots, \mathcal{O}(\Pi_{PK_D,SK_{D-1},K_D}) \right)_{i=0,\dots,D-1} \approx_c \left( \mathcal{O}(\Pi'_{PK_1,K_1}), \dots, \mathcal{O}(\Pi'_{PK_D,K_D}) \right)_{i=0,\dots,D-1} \quad (1)$$

We will do this in a series of hybrids starting with replacing  $\mathcal{O}(\Pi_{PK_D,SK_{D-1},K_D})$  in the LHS by  $\mathcal{O}(\Pi'_{PK_D,K_D})$ , leaving everything else in tact. (Do you see why this is the right order to do the hybrids?)

OK, so how does one show that given  $SK_{D-1}$ , an obfuscation of  $\mathcal{O}(\Pi_{PK_D,SK_{D-1},K_D})$  is indistinguishable from an obfuscation of  $\mathcal{O}(\Pi'_{PK_D,K_D})$ ? We again go through hybrids and a puncturing argument, this time

as many hybrids as the number of possible inputs to either circuit. Note that this is exponentially many! Let's number the inputs from 0 to  $\ell := 2^{|C|+1} - 1$  where  $|C|$  is the bit-length of a ciphertext.

*Hybrid 0.* Define  $\Pi^{(0)} := \Pi_{\text{PK}_D, \text{SK}_{D-1}, K_D}$ . The adversary sees  $\text{SK}_{D-1}$  and  $\mathcal{O}(\Pi^{(0)})$ .

...

*Hybrid  $i$ .* Define  $\Pi^{(i)}$  to be the program that on input  $(C_0, C_1)$ , checks if it is less than  $i$ . If yes, output  $\Pi'_{\text{PK}_D, K_D}(C_0, C_1)$ , else output  $\Pi_{\text{PK}_D, \text{SK}_{D-1}, K_D}(C_0, C_1)$ . The adversary sees  $\text{SK}_{D-1}$  and  $\mathcal{O}(\Pi^{(i)})$ .

...

*Hybrid  $\ell + 1$ .* Define  $\Pi^{(\ell+1)}$  to be the program that on input  $(C_0, C_1)$ , checks if it is less than  $\ell + 1$ . If yes, output  $\Pi'_{\text{PK}_D, K_D}(C_0, C_1)$ , else output  $\Pi_{\text{PK}_D, \text{SK}_{D-1}, K_D}(C_0, C_1)$ . The adversary sees  $\text{SK}_{D-1}$  and  $\mathcal{O}(\Pi^{(\ell+1)})$ .

Note that in the first hybrid, the adversary sees an obfuscation of  $\Pi_{\text{PK}_D, \text{SK}_{D-1}, K_D}$  and in the last, she sees an obfuscation of  $\Pi'_{\text{PK}_D, K_D}$ . If an adversary can distinguish between the first and last hybrids with advantage  $\epsilon$ , she must be able to distinguish between some two adjacent hybrids with advantage at least  $\epsilon/2^\ell$ . We will (try to) show that this is impossible invoking the security of (a) the IO scheme, (b) the encryption scheme and (c) the PRF.

*The Problem.* The first step is to replace the PRF key  $K$  with a key punctured at input  $i$ ; this is OK by IO. The next step is to replace the value of the PRF at  $i$  with a truly random string; this is OK by the puncturing security of the PRF. Now, the output of the programs on input  $i$  is either

$$\text{Enc}_{\text{PK}_i}(\text{Dec}_{\text{SK}_{i-1}}(C_0) \text{ NAND } \text{Dec}_{\text{SK}_{i-1}}(C_1); r) \quad \text{or} \quad \text{Enc}_{\text{PK}_i}(0; r),$$

in the two circuits respectively. In the last step, we want to switch from the first to the second of these ciphertexts using the semantic security of the encryption scheme.

In each of these steps, one uses the subexponential security of the underlying primitive: IO, PRF, or the encryption scheme. The first two are OK by increasing the key sizes and the security parameters sufficiently, but we run into trouble with the encryption scheme. Recall that we will get a distinguisher for the encryption scheme with advantage  $1/(\text{poly}(\lambda) \cdot 2^{\text{ct}+1})$  where  $\text{ct}$  is a ciphertext in the scheme. However, *any encryption scheme can be broken with such small advantage!* (Do you see how?)

*The Fix.* We will rely on an additional primitive: a perfectly lossy encryption scheme. This is an encryption scheme with two types of keys: regular keys  $\text{PK}$  that come with a matching secret decryption key  $\text{SK}$ ; and lossy keys  $\widetilde{\text{PK}}$  where a random ciphertext of 0 is *identically distributed* to a random ciphertext of 1. Moreover, a random  $\text{PK}$  and a random  $\widetilde{\text{PK}}$  are computationally indistinguishable.

Armed with a lossy encryption scheme, we proceed as follows. We change the big hybrids (see equation 1) so that after every step  $i$ , we change the public keys from real to lossy. That is, while before, we were trying to go from an obfuscation of  $\mathcal{O}(\Pi_{\text{PK}_D, \text{SK}_{D-1}, K_D})$  to an obfuscation of  $\mathcal{O}(\Pi'_{\text{PK}_D, K_D})$  directly, we insert two hybrids in between, namely  $\mathcal{O}(\Pi_{\widetilde{\text{PK}}_D, \text{SK}_{D-1}, K_D})$  and  $\mathcal{O}(\Pi'_{\widetilde{\text{PK}}_D, K_D})$ . We then proceed as follows.

$\mathcal{O}(\Pi_{\text{PK}_D, \text{SK}_{D-1}, K_D}) \approx \mathcal{O}(\Pi_{\widetilde{\text{PK}}_D, \text{SK}_{D-1}, K_D})$ . This follows from the indistinguishability of lossy and real public keys.

$\mathcal{O}(\Pi'_{\text{PK}_D, K_D}) \approx \mathcal{O}(\Pi'_{\widetilde{\text{PK}}_D, K_D})$ . This follows from the indistinguishability of lossy and real public keys.

$\mathcal{O}(\Pi_{\widetilde{\text{PK}}_D, \text{SK}_{D-1}, K_D}) \approx \mathcal{O}(\Pi'_{\widetilde{\text{PK}}_D, K_D})$ . Here, we follow the hybrid argument as above except we note that in the step where one invokes the IND-CPA security of encryption, we have *perfect* security! Thus, a distinguisher breaking the indistinguishability of the two corresponding hybrids with *any non-zero* advantage gives us a contradiction.

This observation finishes the construction and the proof.

**Remark** One might wonder if the lossy encryption scheme in the construction is necessary; turns out it is. Bitansky, Degwekar and Vaikuntanathan constructed an oracle world where IO and OWF (even OWP) exist, but where  $\text{SZK} = \text{P}$ . Since every FHE scheme can be broken in  $\text{SZK}$ , it turns out that in this oracle world, FHE does not exist. Therefore, there are no *black-box* FHE constructions from IO and OWP. (The notion of black-boxness has to be formulated carefully to allow obfuscating programs that invoke the OWF/OWP; for more details, we refer the reader to the paper of Bitansky, Degwekar and Vaikuntanathan.)

## 4 IO Engineering II: Using (Zero Knowledge) Proofs

### 4.1 Bootstrapping IO: From Shallow Circuits to All Circuits

There are many ways to construct an IO for all poly-size circuits starting from an IO for relatively simple complexity classes. We start by describing perhaps the simplest of these, which also gives us more practice on how to use IO.

**From Fully Homomorphic Encryption.** We can show that if there is IO for  $\text{NC}^1$  circuits and an FHE scheme whose decryption can be done in  $\text{NC}^1$ , then there is an IO for all of  $P$ .

The idea is to encrypt the circuit  $C$  using an FHE key into  $\widehat{C} = \text{FHE.Enc}(\text{sk}, C)$  so that on any input  $x$ , one can homomorphically compute  $\widehat{C(x)} \in \text{FHE.Enc}(\text{sk}, C(x))$ . The question is how to recover  $C(x)$ . For that purpose, the first, and somewhat naïve, idea is to provide an obfuscation of the following program  $P_{\text{sk}}$ .

On input  $y$ : Output  $\text{FHE.Dec}(\text{sk}, y)$ .

Clearly, this is insecure as  $\mathcal{O}(P_{\text{sk}})$  can be used to decrypt  $\widehat{C}$  and recover the original circuit  $C$ .

A better idea is to obfuscate the “check-then-decrypt” program  $P_{\text{sk},c}$  which has the secret key  $\text{sk}$  and the ciphertext  $c = \widehat{C}$  encoded, and takes a ciphertext  $y$ , an input  $x$ , and a “proof”  $\pi$  as input. If the proof checks out, it decrypts  $y$ . Importantly for us, the proof will simply be the transcript of the homomorphic computation that takes  $\widehat{C}$  to  $\widehat{C(x)}$ , and verifying it can be done in parallel, i.e. in  $\text{NC}^1$ . So:

On input  $(y, \pi)$ : Check the proof  $\pi$ ; If it fails, output  $\perp$ , else output  $\text{FHE.Dec}(\text{sk}, y)$ .

This “should” work but it is unclear how to prove it secure using the IO guarantee.

Instead, we resort to the two-track trick ala Naor-Yung. In other words, the obfuscation of  $C$  is a pair of ciphertexts

$$c_0 = \text{FHE.Enc}(\text{pk}_0, C) \text{ and } c_1 = \text{FHE.Enc}(\text{pk}_1, C)$$

and a program that, on input  $c'_0, c'_1, x, \pi$  as input, checks the proof, and if the proof verifies, decrypts  $c'_0$  using  $\text{sk}_0$ .

This turns out to do the trick. To see why, we construct a sequence of hybrids.

*Hybrid 1.*  $c_0$  and  $c_1$  as above and  $\mathcal{O}(P_{\text{sk}_0, c_0, c_1})$ .

*Hybrid 2.* Change  $c_1$  to be an encryption of  $C_1$  under  $\text{sk}_1$ . OK by semantic security since  $\text{sk}_1$  is never used in these hybrids.

*Hybrid 3.* Change the program to be an obfuscation of  $\mathcal{O}(P_{\text{sk}_1, c_0, c_1})$  which decrypts using  $\text{sk}_1$ . These two programs are functionally equivalent by the perfect correctness of FHE evaluation and the functional equivalence of  $C_0$  and  $C_1$ . Thus, IO guarantees that hybrids 2 and 3 are computationally indistinguishable.

*Hybrid 4.* Change  $c_0$  to be an encryption of  $C_1$  under  $\text{sk}_0$ . OK by semantic security since  $\text{sk}_0$  is never used in these hybrids.

*Hybrid 5.* Once again, change the program to be an obfuscation of  $\mathcal{O}(P_{\text{sk}_0, c_0, c_1})$  which decrypts using  $\text{sk}_0$ . The distribution in hybrid 5 is that of obfuscating program  $C_1$ . This finishes the proof.

## 4.2 (Simulation-Sound) NIZK Proofs from NIZK Proofs

**Achieving Statistical Simulation-Soundness.** We use the following standard compiler from any NIZK into one that is statistically simulation-sound. The CRS consists of the original CRS (of an underlying NIZK system) together with a commitment to 0, i.e.  $c = \text{Com}(0^n; r)$  where  $0^n$  is a special string that cannot be a plausible NP statement. The prover, given  $(x, w)$  shows that either  $x \in L$  or  $\exists r$  such that  $c = \text{Com}(x; r)$ . Soundness follows from that of the underlying NIZK together with the fact that  $c$  is not a commitment of any plausible NP statement. The zero knowledge simulator sets  $c$  to be a commitment of  $x^*$  and uses the corresponding randomness to generate a simulated proof. Even given such a proof, there do not exist accepting proofs of any  $x' \notin R_L$ ,  $x' \neq x^*$ , by the statistical soundness of the underlying NIZK system and the perfect binding of the commitment scheme.

NIZK proofs in turn can be constructed from several algebraic assumptions including quadratic residuosity, LWE, factoring, decisional Diffie-Hellman and more.

## 4.3 Functional Encryption from IO

First, we define functional encryption. Informally, functional encryption allows the release of “partial” decryption keys, which allows one to reveal specified functions  $f$  of the underlying plain-text.

**Definition 3** (Functional Encryption). *A functional encryption scheme is given by a tuple of PPT machines (Setup, KeyGen, Enc, Dec) satisfying the following conditions:*

**Correctness.** *For all polynomial-size circuits  $f$  and inputs  $x$ ,*

$$\Pr [\text{Dec}(\text{sk}_f, \text{Enc}(\text{mpk}, x)) = f(x)] = 1,$$

*where  $(\text{mpk}, \text{msk}) \leftarrow \text{Setup}(1^\lambda)$  and  $\text{sk}_f \leftarrow \text{KeyGen}(\text{msk}, f)$ .*

**Indistinguishable Security (Adaptive).** *For every (stateful) PPT adversary  $\mathcal{A}$ , there is a negligible function  $\mu(\cdot)$  and polynomially-bounded  $m_1(\cdot)$  and  $m_2(\cdot)$  such that the adversary has advantage at most  $\mu(\lambda)$  in the experiment  $\{\text{Exp}_b^{\mathcal{A}}(1^\lambda)\}_{b \in \{0,1\}}$  defined as follows:*

1. *The challenger runs  $(\text{mpk}, \text{msk}) \leftarrow \text{Setup}(1^\lambda)$  and sends  $\text{mpk}$  to  $\mathcal{A}$ .*
2. *Adversary  $\mathcal{A}$  sequentially and adaptively chooses circuits  $f_i$  to send to the challenger, and the challenger sends back  $\text{sk}_{f_i} \leftarrow \text{KeyGen}(\text{msk}, f_i)$  for each  $i \in [m_1(\lambda)]$ .*
3. *Adversary  $\mathcal{A}$  sends two inputs  $x_0, x_1$  to the challenger such that  $f_i(x_0) = f_i(x_1)$  for all  $i \in [m_1(\lambda)]$ .*
4. *The challenger samples  $c \leftarrow \text{Enc}(\text{mpk}, x_b)$  and sends  $c$  to  $\mathcal{A}$ .*



5. Adversary  $\mathcal{A}$  sequentially and adaptively chooses more circuits  $g_i$  to send to the challenger, and the challenger sends back  $sk_{g_i} \leftarrow \text{KeyGen}(msk, g_i)$  for all  $i \in [m_2(\lambda)]$ , as long as  $g_i(x_0) = g_i(x_1)$  for all  $i \in [m_2(\lambda)]$ .
6. Adversary  $\mathcal{A}$  concludes by sending some bit  $b' \in \{0, 1\}$  to the challenger. This bit  $b'$  is defined as the output of the experiment.

We define the advantage of  $\mathcal{A}$  to be

$$\left| \Pr [\text{Exp}_0^{\mathcal{A}}(1^\lambda) = 1] - \Pr [\text{Exp}_1^{\mathcal{A}}(1^\lambda) = 1] \right|,$$

which we require to be at most  $\mu(\lambda)$ .

One may define a weaker notion of security called selective (indistinguishable) security, where the adversary's challenge inputs are declared at the very beginning of the experiment (see e.g. [?] in comparing selective and adaptive security). Also note that a "master decryption key" can be released by obtaining the decryption key for the identity function.

A first idea for a construction of functional encryption from IO is to let the functional key of a circuit  $f$  (represented as a string) be a digital signature of  $f$  using the master secret key  $msk$  as the signing key. An encryption of a message  $m$  is the obfuscation of the program  $P_{m,mpk}$  that takes as input  $(f, \sigma)$  and checks if  $\sigma$  is a valid signature of  $f$  under the verification key  $mpk$ . If yes, the program outputs  $f(m)$ ; otherwise, the program outputs  $\perp$ . While this is a natural idea, it is unclear how to make it work using IO. The reason, roughly speaking, is that IO is useful when one can perhaps "puncture the verification key" in some way to ensure that some  $f$ 's do not have any valid signature.

**Open Problem 4.** *Can the above idea go through, perhaps by constructing an appropriate puncturable signature scheme?*

While one may be able to operationalize the above idea, we take a different route here. Assume the existence of statistically simulation-sound NIZK proof system and a CPA encryption scheme. The construction works as follows.

1. The setup algorithm generates two public keys  $PK_0, PK_1$  and the matching private keys  $SK_0, SK_1$  of a CPA-secure encryption scheme, as well as a common (or structured) random string  $CRS$  of a statistically simulation-sound NIZK proof system. The setup is then

$$MPK = (PK_0, PK_1, CRS) \text{ and } MSK = (SK_0, SK_1).$$

2. The encryption algorithm, given  $x$ , encrypts  $x$  with both public keys and attaches a NIZK proof that the two ciphertexts encrypt the same message.

$$C_b \leftarrow \text{Enc}(PK_b, x; r_b) \text{ and } \pi \leftarrow P(CRS, C_0, C_1, x, r_0, r_1)$$

3. The key generation algorithm, given  $MSK$  and a circuit  $f$ , writes down a program  $\Pi_{CRS, SK_0, f}$  that takes as input  $(C_0, C_1, \pi)$ .

$\Pi_{\text{CRS}, \text{SK}_0, f}(C_0, C_1, \pi)$  **does the following:**

- Verify  $\pi$  as a proof that  $C_0$  and  $C_1$  encrypt the same message (using CRS).
- If valid, output  $f(\text{Dec}(\text{SK}_0, C_0))$ .
- Else, output  $\perp$ .

We then define

$$\text{SK}_f = \mathcal{O}(\Pi_{\text{CRS}, \text{SK}_0, f}).$$

4. The decryption algorithm simply feeds the ciphertext into the secret key  $\text{SK}_f$  which is the obfuscated program.

The (selective security) proof proceeds via the following sequence of hybrids.

*Hybrid 0.* The adversary declares a challenge  $(x^*, y^*)$ . Construct  $C_0^*$  and  $C_1^*$  by encrypting  $x^*$  under  $\text{PK}_0$  and  $\text{PK}_1$  respectively. Let CRS be the real CRS. Produce a real proof  $\pi$ .

*Hybrid 1.* Replace the common reference string with the simulated  $\text{CRS}^*$ . Produce a simulated proof  $\pi^*$ . This is indistinguishable by the ZK property of the proof system.

*Hybrid 2.* Replace  $C_1^*$  to be the encryption of  $y^*$ . The proof is still simulated. This is indistinguishable by CPA-security since  $\text{SK}_1$  is never used.

*Hybrid 3.* Replace the functional keys with ones that decrypt using  $\text{SK}_1$ . This is OK because of the IO guarantee and because the circuits that decrypt using  $\text{SK}_0$  and  $\text{SK}_1$  are identical in functionality. Here, we crucially use the fact that the NIZK proof system is statistically simulation-sound, and that  $f(x^*) = f(y^*)$  for every function  $f$  for which the adversary obtains the functional key.

*Hybrid 4.* Replace  $C_0^*$  to be the encryption of  $y^*$ . The proof is still simulated. This is indistinguishable by CPA-security since  $\text{SK}_0$  is never used.

*Hybrid 5.* Replace the functional keys with ones that decrypt using  $\text{SK}_0$ . This is OK because of the same argument as in Hybrid 2 vs 3.

*Hybrid 6.* Replace the common reference string with the real CRS and the proof with the real proof. This is indistinguishable by the ZK property of the proof system.

Hybrid 6 is the real game where  $y^*$  is encrypted instead of  $x^*$  in Hybrid 0. Since hybrids 0 and 6 are indistinguishable, we are done.