

# PRFs and the GGM construction

Noah Stephens-Davidowitz

June 8, 2023

## 1 Pseudorandom functions

In the last lecture, we saw how to build *stateful* semantically secure secret-key encryption using a PRG. But, we still have not built a true semantically secure encryption scheme under our original definition.

Now, we fix this by showing how to construct an efficiently computable *pseudorandom* function, a PRF. Formally, it is easier to define “a pseudorandom function” as a *keyed* function, also known as a function *family*,  $F_{\mathbf{s}}(\mathbf{r})$ , where  $\mathbf{s} \in \{0, 1\}^*$  and the length of the input  $\mathbf{r}$  depends on the length of  $\mathbf{s}$ . This notation is formally equivalent to writing  $F(\mathbf{s}, \mathbf{r})$ , but we write  $F_{\mathbf{s}}(\mathbf{r})$  to emphasize that we think of  $\mathbf{s}$  as the key and  $\mathbf{r}$  as the input. When the key  $\mathbf{s}$  is  $n$  bits long, our PRF will take  $n$  bits as input and output  $\ell(n)$  bits. (One can be slightly more general and take  $m(n)$  bits of input for any  $m(n) = \text{poly}(n)$ , but it does not change much.)

The security of a PRF is defined in terms of the following two games. In one game, the adversary  $\mathcal{A}$  is given *oracle* access to  $F_{\mathbf{s}}$ . I.e., the challenger samples  $\mathbf{s} \sim \{0, 1\}^n$ , and  $\mathcal{A}$  can make as many queries as it likes (though the number of queries must of course be polynomially bounded, because  $\mathcal{A}$  must run in polynomial time) to  $F_{\mathbf{s}}$ . Each query consists of an  $n$ -bit string  $\mathbf{x}$ , and  $\mathcal{A}$  receives in response an  $F_{\mathbf{s}}(\mathbf{x})$ . Eventually  $\mathcal{A}$  outputs either zero or one.

Super succinct (and convenient!) notation for this game is given by  $b' \leftarrow \mathcal{A}^{F_{\mathbf{s}}}(1^n)$  where  $\mathbf{s} \sim \{0, 1\}^n$ , which means “ $b'$  is the output of  $\mathcal{A}$  on input  $1^n$  with oracle access to  $F_{\mathbf{s}}$  for  $\mathbf{s} \sim \{0, 1\}^n$ .” Giving  $\mathcal{A}$  “oracle access” to  $F_{\mathbf{s}}$  means that  $\mathcal{A}^{F_{\mathbf{s}}}$  is an algorithm that has access to an extra procedure that computes  $F_{\mathbf{s}}$ . The full game is shown below.

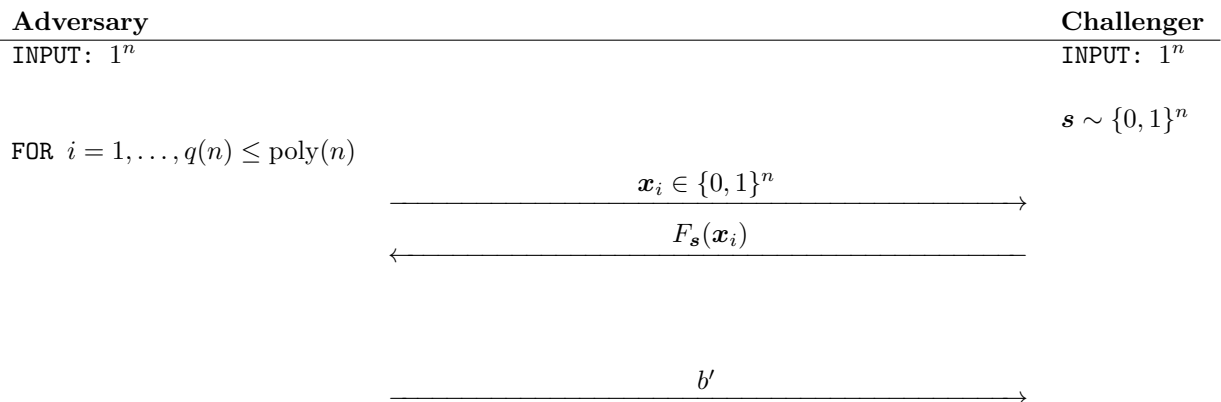


Figure 1: The game in which  $\mathcal{A}$  is given oracle access to  $F_{\mathbf{s}}$  for  $\mathbf{s} \sim \{0, 1\}^n$ .

In the second game, the pseudorandom bits  $F_s(\mathbf{x}_i)$  are replaced by truly random bits  $\mathbf{y}_i \sim \{0, 1\}^{\ell(n)}$ . Specifically, in response to each *distinct* query  $\mathbf{x}$ , the challenger responds with a uniformly random string  $\mathbf{y} \sim \{0, 1\}^{\ell(n)}$ . However, notice that in the game with  $F_s$ , if  $\mathcal{A}$  makes the same query  $\mathbf{x}$  twice, it will get the same answer  $F_s(\mathbf{x})$  both times. So, the challenger does the same thing in our new game. Specifically, it records the queries and makes sure to respond to duplicate queries in a consistent manner. (This makes the game itself a little tedious to write out, but the challenger’s code is literally just implementing the procedure “respond randomly unless you have seen this query before, in which case give the same response that you gave previously.”)

Another way to view this is to consider a *random function*  $H : \{0, 1\}^n \rightarrow \{0, 1\}^{\ell(n)}$ . I.e., let  $\mathcal{H}_{n,\ell} := \{H : \{0, 1\}^n \rightarrow \{0, 1\}^{\ell}\}$  be the set of *all* possible functions from  $n$  bits to  $\ell$  bits. Then  $H \sim \mathcal{H}_{n,\ell}$  is a random function. In other words, for every input  $\mathbf{x} \in \{0, 1\}^n$   $H(\mathbf{x})$  is a uniformly random and independent bit string of length  $\ell$ . We can describe our second game succinctly using the notation  $b' \leftarrow \mathcal{A}^H(1^n)$  for uniformly random  $H \sim \mathcal{H}_{n,\ell(n)}$ . Here is the game written out formally.

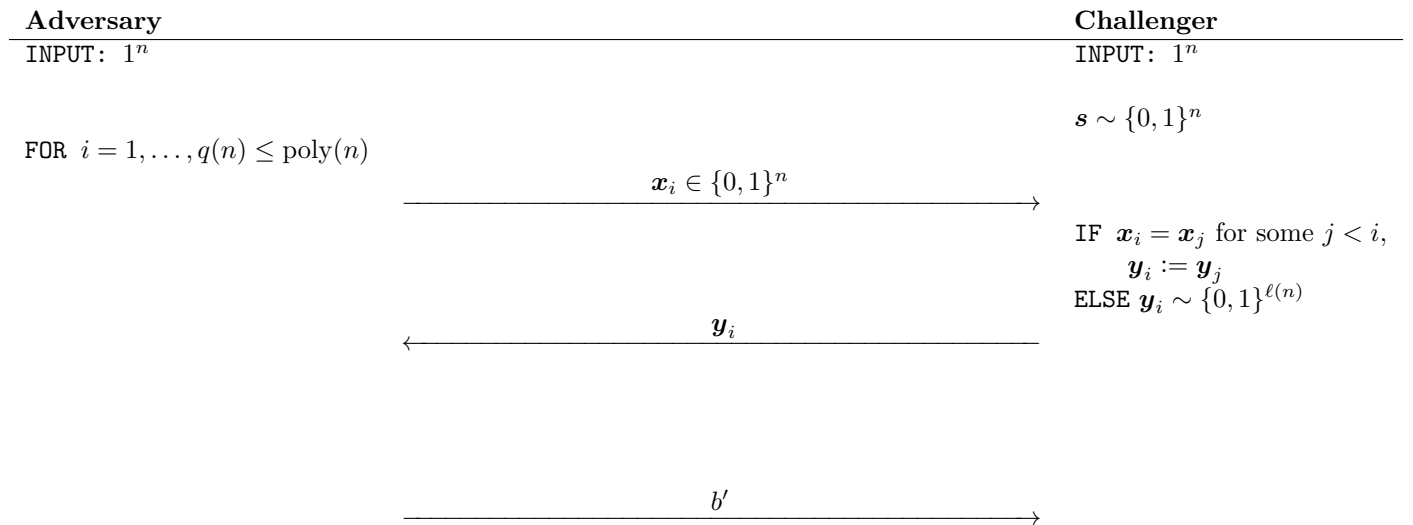


Figure 2: The game in which  $\mathcal{A}$  is given oracle access to  $H \sim \mathcal{H}_{n,\ell(n)}$ .

(You might wonder why I did not simply have the challenger sample  $H \sim \mathcal{H}_{n,\ell(n)}$  and then respond to each query  $\mathbf{x}_i$  with  $H(\mathbf{x}_i)$ . This *would* yield an equivalent game. But, a random function  $H$  is not efficiently computable in general! And, it will cause problems for us later if our challenger does not run in polynomial time. So, we avoid this mess via a “lazy sampling” technique. I.e., rather than choose a fixed random function  $H$  at the beginning of the game and then computing  $H(\mathbf{x})$  repeatedly (which cannot be done efficiently), the challenger “only decides what  $H(\mathbf{x})$  should be when she needs to.”)

With this, we can define a PRF. Notice that the security definition is exactly the same as saying “no PPT adversary has non-negligible advantage in distinguishing the first game above from the second game above.” We have just used this very convenient oracle notation  $\mathcal{A}$  to write that quite succinctly.

**Definition 1.1** (Pseudorandom function). *A function  $F_s : \{0, 1\}^n \rightarrow \{0, 1\}^{\ell(n)}$  for  $\mathbf{s} \in \{0, 1\}^*$  and*

$n := |\mathbf{s}|$  is a pseudorandom function (PRF) if it satisfies the following.

1. **Efficiently computable.** There is a PPT algorithm  $\mathcal{A}$  such that  $\mathcal{A}(\mathbf{s}, \mathbf{x}) = F_{\mathbf{s}}(\mathbf{x})$  for all integers  $n \in \mathbb{N}$ , keys  $\mathbf{s} \in \{0, 1\}^*$ , and inputs  $\mathbf{x} \in \{0, 1\}^{|\mathbf{s}|}$ .
2. **Pseudorandom.** For all PPT algorithms  $\mathcal{A}$ , there exists a negligible function  $\varepsilon(n)$  such that for all  $n \in \mathbb{N}$ ,

$$\left| \Pr_{\mathbf{s} \sim \{0, 1\}^n} [\mathcal{A}^{F_{\mathbf{s}}}(1^n) = 1] - \Pr_{H \sim \mathcal{H}_{n, \ell(n)}} [\mathcal{A}^H(1^n) = 1] \right| \leq \varepsilon(n).$$

## 2 The GGM construction of a PRF from a PRG

We now show how to build a PRF from a PRG. The (elegant) construction that we will use is due to Goldreich, Goldwasser, and Micali [GGM86], and is typically just called the GGM PRF.

To get some intuition for the construction, let's recall our idea from last class for "getting more pseudorandom bits" out of a PRG. I.e., given a PRG  $G$  that maps  $n$  bits to  $2n$  bits, we define  $G_0(\mathbf{s})$  and  $G_1(\mathbf{s})$  as the first  $n$  bits of the output and the second  $n$  bits of the output respectively. Then, to get  $rn$  pseudorandom bits from  $G$  using a seed  $\mathbf{s}_0 \sim \{0, 1\}^n$  of length  $n$ , we define  $\mathbf{s}_i := G_0(\mathbf{s}_{i-1})$  and  $\mathbf{x}_i := G_1(\mathbf{s}_{i-1})$ . Notice that these quantities are efficiently computable for polynomial  $i$ . Our pseudorandom bits are then  $(\mathbf{x}_1, \dots, \mathbf{x}_r) \in \{0, 1\}^{rn}$ .

We can equivalently think of  $G$  and  $\mathbf{s}$  as defining a *function*,  $F_{G, \mathbf{s}}(r)$  that takes as input  $r \in \mathbb{N}$  and outputs  $\mathbf{x}_r$ . Unfortunately, it takes time proportional to  $r$  to compute  $F_{G, \mathbf{s}}(r)$ . We are therefore effectively restricted to computing this function on polynomially many inputs  $r$ . Equivalently, we can think of the domain of  $F_{G, \mathbf{s}}$  as  $\{0, 1\}^m$ , where  $m = O(\log n)$  so that  $|\{0, 1\}^m| = \text{poly}(n)$ .

To get a PRF, we need to handle  $n$ -bit inputs. I.e., instead of taking an  $n$ -bit random seed  $\mathbf{s} \in \{0, 1\}^n$  and expanding it into  $\text{poly}(n)$  pseudorandom bits, a PRF in some sense takes an  $n$ -bit random seed and expands it into  $\ell \cdot 2^n$  pseudorandom bits! And we must have efficient access to these bits.

The problem with  $F_{G, \mathbf{s}}$  is of course that it needs to compute all of the values  $\mathbf{s}_1, \dots, \mathbf{s}_{r-1}$  before computing  $\mathbf{x}_r$ . As a data structure, it is a *linked list*—where the entries in the list are indexed by  $\mathbf{s}_i$  and labeled with  $\mathbf{x}_i$ . To find the label  $\mathbf{x}_r$  of the  $r$ th node in the list, we use  $G_0$  to iteratively compute the indices  $\mathbf{s}_1, \dots, \mathbf{s}_{r-1}$ , before finally using  $G_1$  to compute  $\mathbf{x}_r$ . The GGM PRF replaces this linked list with a balanced binary *tree*. (You can make a career in computer science just by finding places to replace linked lists with trees :).)

Here is the construction. Let's write  $\mathbf{s}_0 := G_0(\mathbf{s})$  and  $\mathbf{s}_1 := G_1(\mathbf{s})$ . (Note that I have redefined things here.) Then,  $\mathbf{s}_{00} := G_0(\mathbf{s}_0)$ ,  $\mathbf{s}_{01} := G_1(\mathbf{s}_0)$ ,  $\mathbf{s}_{10} := G_0(\mathbf{s}_1)$ , and  $\mathbf{s}_{11} := G_1(\mathbf{s}_1)$ . (Here and below, we will use the notation  $\mathbf{rt}$  for the concatenation of the bit strings  $\mathbf{r}$  and  $\mathbf{t}$ .) More generally,  $\mathbf{s}_{r0} := G_0(\mathbf{s}_r)$  and  $\mathbf{s}_{r1} := G_1(\mathbf{s}_r)$ . I.e., unraveling the definition, we have for  $\mathbf{r} := (r_1, r_2, \dots, r_n) \in \{0, 1\}^n$

$$\mathbf{s}_{\mathbf{r}} := G_{r_n}(G_{r_{n-1}}(\dots(G_{r_1}(\mathbf{s}) \dots))).$$

Then, the GGM PRF  $F_{\mathbf{s}} : \{0, 1\}^n \rightarrow \{0, 1\}^n$  is defined as  $F_{\mathbf{s}}(\mathbf{r}) := \mathbf{s}_{\mathbf{r}}$ . (It's quite important that the input size is fixed for fixed  $\mathbf{s}$ ! If we reveal  $\mathbf{s}_{\mathbf{r}}$  for strings  $\mathbf{r}$  with length, say,  $n - 1$  as well, then the adversary can certainly distinguish, say,  $\mathbf{s}_{r0}$  from random, since  $\mathbf{s}_{r0} = G_0(\mathbf{s}_r)$ . So, it is crucial that an adversary making queries to  $F_{\mathbf{s}}$  will only learn  $\mathbf{s}_{\mathbf{r}}$  for strings  $\mathbf{r}$  that have length exactly  $n$ .) See Figure 3.

To prove the security of this construction, we will need the following lemma, which we proved in the previous lecture.

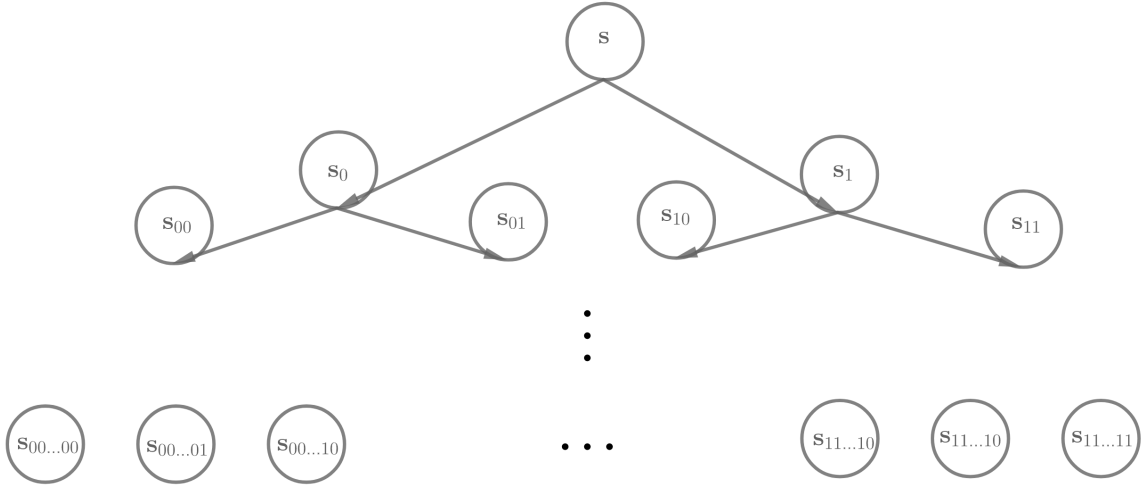


Figure 3: The GGM PRF.

**Lemma 2.1.** *If  $G$  is a length-doubling PRG, then for any PPT  $\mathcal{A}$  and any polynomial  $p(n)$ , there exists a negligible  $\varepsilon(n)$  such that for all  $n \in \mathbb{N}$ ,*

$$\left| \Pr_{\mathbf{s}_1, \dots, \mathbf{s}_{p(n)} \sim \{0,1\}^n} [\mathcal{A}(1^n, G(\mathbf{s}_1), \dots, G(\mathbf{s}_{p(n)})) = 1] - \Pr_{\mathbf{x}_1, \dots, \mathbf{x}_{p(n)} \sim \{0,1\}^{2n}} [\mathcal{A}(1^n, \mathbf{x}_1, \dots, \mathbf{x}_{p(n)}) = 1] \right| \leq \varepsilon(n).$$

**Theorem 2.2.** *If  $G$  is a PRG, then  $F_s$  as defined above is a PRF.*

*Proof.* The proof is by a hybrid argument. Let  $F_0$  be the oracle  $F_s$  sampled honestly. I.e.,  $F_0$  has a key  $\mathbf{s} \sim \{0,1\}^n$  sampled uniformly at random, and  $F_0$  then returns  $F_s(\mathbf{r})$  for any query  $\mathbf{r}$  made by the adversary  $\mathcal{A}$ .  $F_1$  has two keys  $\mathbf{s}_0, \mathbf{s}_1 \sim \{0,1\}^n$  sampled independently and uniformly at random, and it behaves like the GGM PRF with the pseudorandom strings in the first level below the root replaced by the pseudorandom  $\mathbf{s}_0$  and  $\mathbf{s}_1$ . I.e., we “remove the root and split the tree into two trees,” as in Figure 4.

More generally  $F_i$  has  $2^i$  keys, given by  $\mathbf{s}_{00\dots 0}, \dots, \mathbf{s}_{11\dots 1} \sim \{0,1\}^n$  all sampled uniformly and independently at random, and to compute  $F_i$ , we compute the GGM PRF, except that we replace the whole  $i$ th level of the tree by these random keys, rather than the pseudorandom elements in the “honest” construction. Notice in particular that  $F_n$  is a purely random oracle.

So, if some adversary  $\mathcal{A}$  has non-negligible advantage  $\varepsilon$  in breaking the PRF (i.e., if it distinguishes the PRF from a random oracle with probability  $\varepsilon$ ), then there must be some index  $1 \leq i \leq n$  such that  $\mathcal{A}$  distinguishes  $F_i$  from  $F_{i-1}$  with probability at least  $\varepsilon/n$ .

It therefore suffices to show how to use an adversary  $\mathcal{A}$  that distinguishes between  $F_i$  and  $F_{i-1}$  in order to build an adversary  $\mathcal{B}$  that breaks the security of the PRG. Or, rather, we will build an adversary  $\mathcal{B}$  that breaks the game described in Lemma 2.1. To do this, we must decide what  $p$  should be.

Intuitively, we want to use the  $p$  strings  $\mathbf{x}_1, \dots, \mathbf{x}_p$  that  $\mathcal{B}$  receives as input to replace the strings labeling the nodes in the  $i$ th row of Figure 3. Then, if the  $\mathbf{x}_j$  are random, we will faithfully reproduce  $F_i$ . If they are pseudorandom, we will faithfully reproduce  $F_{i-1}$ . So, an  $\mathcal{A}$  that distinguishes between an  $F_i$  and  $F_{i-1}$  will yield a  $\mathcal{B}$  that breaks our PRG.

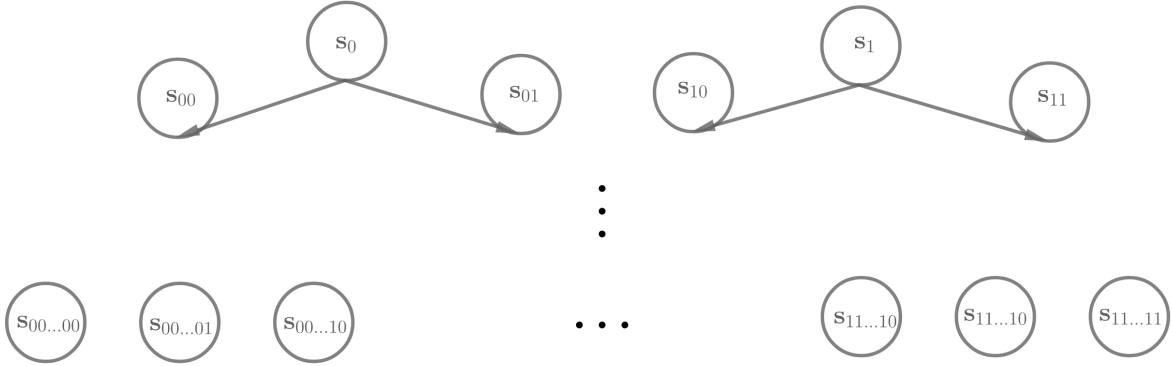


Figure 4: The GGM tree with the root removed, used in Game 1. Here,  $\mathbf{s}_0$  and  $\mathbf{s}_1$  are both sampled uniformly at random

The only problem is that the number of nodes in the  $i$ th row is  $2^i$ , which is superpolynomial for large  $i$ . We need  $p$  to be polynomial in  $n$ . (This isn't just a technical issue. Lemma 2.1 is simply false for superpolynomial  $p$ .) To handle this, we will use the lazy sampling idea that we discussed above. I.e., we will not assign a string to each node in the  $i$ th row right away. Instead, we will only assign strings to nodes that are queried, and we will do so only when they are queried.

So, let  $p$  be a polynomial that bounds the number of queries that  $\mathcal{A}$  makes to its oracle.  $\mathcal{B}$  receives as input the strings  $\mathbf{x}_1 := (\mathbf{x}_{1,0}, \mathbf{x}_{1,1}), \dots, \mathbf{x}_p := (\mathbf{x}_{p,0}, \mathbf{x}_{p,1}) \in \{0, 1\}^{2n}$ . (To be clear, we have simply split the  $2n$ -bit strings into pairs of  $n$ -bit strings for convenience.) When  $\mathcal{B}$  receives its first oracle query  $\mathbf{r} \in \{0, 1\}^n$  from  $\mathcal{A}$ , it behaves as follows. Let  $\mathbf{r}' \in \{0, 1\}^{i-1}$  be the first  $i-1$  bits of  $\mathbf{r}$ .  $\mathcal{B}$  sets  $\mathbf{s}_{\mathbf{r}'0} := \mathbf{x}_{1,0}$  and  $\mathbf{s}_{\mathbf{r}'1} := \mathbf{x}_{1,1}$ . It then computes  $\mathbf{s}_{\mathbf{r}} := G_{r_n}(G_{r_{n-1}}(\dots G_{r_{i+1}}(\mathbf{s}_{\mathbf{r}'r_i}) \dots))$ , just like  $F_i$ .

For the  $j$ th oracle query  $\mathbf{r}$ , we again define  $\mathbf{r}' \in \{0, 1\}^{i-1}$  to be the first  $i-1$  bits of  $\mathbf{r}$ . If  $\mathbf{s}_{\mathbf{r}'0}, \mathbf{s}_{\mathbf{r}'1}$  have not been defined previously, then  $\mathcal{B}$  sets  $\mathbf{s}_{\mathbf{r}'0} := \mathbf{x}_{j,0}$  and  $\mathbf{s}_{\mathbf{r}'1} := \mathbf{x}_{j,1}$ . It then computes  $\mathbf{s}_{\mathbf{r}} := G_{r_n}(G_{r_{n-1}}(\dots G_{r_{i+1}}(\mathbf{s}_{\mathbf{r}'r_i}) \dots))$  as before.

After at most  $p$  queries  $\mathcal{A}$  will return either 0 or 1, and  $\mathcal{B}$  simply outputs the same.

Clearly,  $\mathcal{B}$  is efficient. To show that the advantage of  $\mathcal{B}$  in the game defined in Lemma 2.1 is equal to the advantage of  $\mathcal{A}$  in the  $i$ th hybrid, we simply need to observe that (1) when  $(\mathbf{x}_{j,0}, \mathbf{x}_{j,1}) = G(\mathbf{s}_j)$  for independent  $\mathbf{s}_j$ , then the responses of  $\mathcal{B}$  to the queries made by  $\mathcal{A}$  are distributed identically to the responses of  $F_{i-1}$ ; and (2) when the  $\mathbf{x}_{j,b}$  are uniformly and independently random, then we need to show that the responses are distributed identically to the responses of  $F_i$ . Both observations follow immediately from the definitions of  $F_i$  and the game in Lemma 2.1.  $\square$

## References

- [GGM86] Oded Goldreich, Shafi Goldwasser, and Silvio Micali. How to construct random functions. *J. ACM*, 33(4), 1986.