

# Aggregate Pseudorandom Functions and Connections to Learning

Aloni Cohen<sup>1</sup>, Shafi Goldwasser<sup>1,2</sup>, and Vinod Vaikuntanathan<sup>1</sup>

<sup>1</sup> MIT

<sup>2</sup> Weizmann Institute of Science

**Abstract.** In the first part of this work, we introduce a new type of pseudo-random function for which “aggregate queries” over exponential-sized sets can be efficiently answered. An example of an aggregate query may be the product of all function values belonging to an exponential-sized interval, or the sum of all function values on points for which a polynomial time predicate holds. We show how to use algebraic properties of underlying classical pseudo random functions, to construct aggregatable pseudo random functions for a number of classes of aggregation queries under cryptographic hardness assumptions. On the flip side, we show that certain aggregate queries are impossible to support.

In the second part of this work, we show how various extensions of pseudo-random functions considered recently in the cryptographic literature, yield impossibility results for various extensions of machine learning models, continuing a line of investigation originated by Valiant and Kearns in the 1980s and 1990s. The extended pseudo-random functions we address include constrained pseudo random functions, aggregate pseudo random functions, and pseudo random functions secure under related-key attacks.

## 1 Introduction

Pseudo-random functions (PRF), introduced by Goldreich, Goldwasser and Micali [GGM86], are a family of indexed functions for which there exists a polynomial-time algorithm that, given an index (which can be viewed as a secret key) for a function, can evaluate it, but no probabilistic polynomial-time algorithm without the secret key can distinguish the function from a truly random function – even if allowed oracle query access to the function. Pseudo-random functions have been shown over the years to be useful for numerous cryptographic applications. Interestingly, aside from their cryptographic applications, PRFs have also been used to show impossibility of computational learning in the membership queries model [Val84], and served as the underpinning of the proof of Razborov and Rudich [RR97] that natural proofs would not suffice for unrestricted circuit lower bounds.

Since their inception in the mid eighties, various *augmented* pseudo random functions with *extra* properties have been proposed, enabling more sophisticated forms of access to PRFs and more structured forms of PRFs. In this

vein, we mention the work of [GGN10] on constructing “huge random objects” which designs PRFs that are guaranteed to maintain some global combinatorial property; the recent works on constrained PRFs<sup>3</sup> [KPTZ13a,BGI14a,BW13a] which can release auxiliary secret keys whose knowledge enables computing the PRF in a restricted number of locations without compromising pseudo-randomness elsewhere; the construction of a PRF [BLMR13] family which is homomorphic with respect to its key; and the construction of related key secure PRFs [BC10,ABPP14]. These constructions yield fundamental objects with often surprising applications to cryptography and elsewhere. A case in point is the truly surprising use of constrained PRFs [SW14], to show that indistinguishability obfuscation can be used to resolve a long-standing problem of deniable encryption, among many others.

In the first part of this paper, we introduce a new type of augmented PRF which we call *aggregate pseudo random functions* (AGG-PRF). An AGG-PRF is a family of indexed functions each associated with a secret key, such that *given the secret key*, one can compute aggregates of the values of the function *over super-polynomially large sets in polynomial time*; and yet without the secret key, access to such aggregated values cannot enable a polynomial time adversary (distinguisher) to distinguish the function from random, even when the adversary can make *aggregate queries*. Note that the distinguisher can request and receive an aggregate of the function values over sets (of possibly super-polynomial size) that she can specify. Examples of aggregate queries can be the sum/product of all function values belonging to an exponential-sized interval, or more generally, the sum/product of all function values on points for which some polynomial time predicate holds. Since the sets over which our function values are aggregated are super-polynomial in size, they cannot be directly computed by simply querying the function on individual points.

We show AGG-PRFs under various cryptographic hardness assumptions (one-way functions and DDH) for a number of types of aggregation operators such as sums and products and for a number of set systems including intervals, hypercubes, and (the supports of) restricted computational models such as decision trees and read-once Boolean formulas. We also show negative results: there are no AGG-PRFs for more expressive set systems such as (the supports of) CNF formulas. For a detailed description of our results, see Section 1.1.

In the second part of this paper, we embark on a study of the connection between the new augmented PRF constructions of recent years (constrained, related-key, aggregate) and the theory of computational learning. We recall at the outset that the fields of cryptography and machine learning share a curious historical relationship. The goals are in complete opposition and at the same time the aesthetics of the models, definitions and techniques bear a striking similarity. For example, a cryptanalyst can attack a cryptosystem using a range of powers from only seeing ciphertext examples to requesting to see decryptions of ciphertexts of her choice. Analogously, machine learning allows different powers to the learner such as random examples versus membership queries and shows

---

<sup>3</sup> Constrained PRFs are also known as Functional PRFs and as Delegatable PRFs.

that certain powers allow learners to learn concepts in polynomial time whereas others will fail. Even more directly, problems which pose challenges for machine learning such as Learning Parity with Noise (LPN) have been used as the underpinning for building secure cryptosystems, and as mentioned above [Val84] observes that the existence of PRFs in a complexity class  $\mathcal{C}$  implies the existence of concept classes in  $\mathcal{C}$  which can not be learned under membership queries, and [KV94] extends this direction to some public key constructions.

In the decades since the introduction of PAC learning, new computational learning models have been proposed, such as the recent “restriction access” model [DRWY12] which allows the learner to interact with the target concept by asking membership queries, but also to obtain an entire circuit that computes the concept on a random subset of the inputs. For example, in one shot, the learner can obtain a circuit that computes the concept class on all  $n$ -bit inputs that start with  $n/2$  zeros. At the same time, the cryptographic research landscape has been swiftly moving in the direction of augmenting traditional PRFs and other cryptographic primitives to include higher functionalities. This brings to mind natural questions:

- *Can one leverage augmented pseudo-random function constructions to establish limits on what can and cannot be learned in augmented machine learning models?*
- *Going even further afield, can augmented cryptographic constructs suggest interesting learning models?*

We address these questions in the second part of this paper. For a detailed description of our findings, see Section 1.2.

### 1.1 Our Results: Aggregate Pseudo Random Functions

Aggregate Pseudo Random Functions (AGG-PRF) are indexed families of pseudo-random functions for which a distinguisher (who runs in time polynomial in the security parameter) can request and receive the value of an aggregate (for example, the sum or the product) of the function values over certain large sets and yet cannot distinguish oracle access to the function from oracle access to a truly random function. At the same time, given the function index (in other words, the secret key), one can compute such aggregates over potentially super-polynomial size sets in polynomial time. Such an efficient aggregation algorithm cannot possibly exist for random functions. Thus, this is a PRF family that is very unlike random functions (in the sense of being able to efficiently aggregate over super-polynomial size sets), and yet is computationally indistinguishable from random functions.

To make this notion precise, we need two ingredients. Let  $\mathcal{F} = \{\mathcal{F}_\lambda\}_{\lambda>0}$  where each  $\mathcal{F}_\lambda = \{f_K : \mathcal{D}_\lambda \rightarrow \mathcal{R}_\lambda\}_{K \in \mathcal{K}_\lambda}$  is a collection of functions on a domain  $\mathcal{D}_\lambda$  to a range  $\mathcal{R}_\lambda$ , computable in time  $\text{poly}(\lambda)$ .<sup>4</sup> The first ingredient is a collection of

<sup>4</sup> In this informal exposition, for the sake of brevity, we will sometimes omit the security parameter and refrain from referring to ensembles.

sets (also called a set system)  $\mathcal{S} = \{S \subseteq \mathcal{D}\}$  over which the aggregates can be efficiently computed given the index  $K$  of the function. The second ingredient is an aggregation function  $\Gamma : \mathcal{R}^* \rightarrow \{0, 1\}^*$  which takes as input a tuple of function values  $\{f(x) : x \in S\}$  for some set  $S \in \mathcal{S}$  and outputs the aggregate  $\Gamma(f(x_1), \dots, f(x_{|S|}))$ .

The sets are typically super-polynomially large, but are efficiently recognizable. That is, for each set  $S$ , there is a corresponding  $\text{poly}(\lambda)$ -size circuit  $C_S$  that takes as input an  $x \in \mathcal{D}$  and outputs 1 if and only if  $x \in S$ .<sup>5</sup> Throughout this paper, we will consider relatively simple aggregate functions, namely we will treat the range of the functions as an Abelian group, and will let  $\Gamma$  denote the group operation on its inputs. Note that the input to  $\Gamma$  is super-polynomially large (in the security parameter  $\lambda$ ), making the aggregate computation non-trivial.

This family of functions, equipped with a set system  $\mathcal{S}$  and an aggregation function  $\Gamma$  is called an aggregate PRF family (AGG-PRF) if the following two requirements hold:

1. *Aggregatability*: There exists a polynomial (in the security parameter  $\lambda$ ) time algorithm that given an index  $K$  to the PRF  $f_K \in \mathcal{F}$  and a circuit  $C_S$  that recognizes a set  $S \in \mathcal{S}$ , can compute  $\Gamma$  over the PRF values  $f_K(x)$  for all  $x \in S$ . That is, it can compute

$$AGG_{K,\Gamma}(S) := \Gamma_{x \in S} f_K(x)$$

2. *Pseudorandomness*: No polynomial-time distinguisher which can specify a set  $S \in \mathcal{S}$  as a query and can receive as an answer either  $AGG_{K,\Gamma}(S)$  for a random function  $f_K \in \mathcal{F}$  or  $AGG_{h,\Gamma}(S)$  for a truly random functions  $h$ , can distinguish between the two cases.

We remark that our notion of aggregate PRFs bears some resemblance to the notion of “algebraic PRFs” defined in the work of Benabbas, Gennaro and Vahlis [BGV11]. In a nutshell, there are two main differences. First, algebraic PRFs support efficient aggregation over very specific subsets, whereas our constructions of aggregate PRFs support expressive subset classes, such as subsets recognized by hypercubes, decision trees and read-once Boolean formulas. Secondly, in the security notion for aggregate PRFs, the adversary obtains access to an oracle that computes the function as well as one that computes the aggregate values over super-polynomial size sets, whereas in algebraic PRFs, the adversary is restricted to accessing the function oracle alone. Our constructions from DDH use an algebraic property of the Naor-Reingold PRF in a similar manner as in [BGV11].

We show a number of constructions of AGG-PRF for various set systems under different cryptographic assumptions. We describe our constructions below, starting from the least expressive set system.

---

<sup>5</sup> All the sets we consider are efficiently recognizable, and we use the corresponding circuit as the representation of the set. We occasionally abuse notation and use  $S$  and  $C_S$  interchangeably.

*Interval Sets.* We first construct AGG-PRFs over interval set systems with respect to aggregation functions that compute any group operation. The construction can be based on any (standard) PRF family.

**Theorem 1 (Intervals from one-way functions).** *Assume one-way functions exist. Then, there exists an AGG-PRF family that maps  $\mathbb{Z}_p$  to a group  $G$ , with respect to a collection of sets defined by intervals  $[a, b] \subseteq \mathbb{Z}_p$  and the aggregation function computing the group operation on  $G$ .*

The construction works as follows. Let  $F : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$  be a (standard) pseudo-random function family based on the existence of one-way functions [GGM86, HILL99]. Construct an AGG-PRF family  $G$  supporting efficient computation of group aggregation functions. Define

$$G(k, x) = F(k, x) - F(k, x - 1)$$

To aggregate  $G$ , set

$$\sum_{x \in [a, b]} G(k, x) = F(k, b) - F(k, a - 1)$$

Given  $k$ , this can be efficiently evaluated. This construction can also be easily constrained to allow evaluation on any point or sub-interval of any interval  $[a, b]$ .

*Hypercubes.* As a warmup, we next construct AGG-PRFs over hypercube set systems. Throughout this section, we take  $\mathcal{D}_\lambda = \{0, 1\}^\ell$  for some polynomial  $\ell = \ell(\lambda)$ . A hypercube  $S_{\mathbf{y}}$  is defined by a vector  $\mathbf{y} \in \{0, 1, \star\}^\ell$  as

$$S_{\mathbf{y}} = \{\mathbf{x} \in \{0, 1\}^\ell : \forall i, y_i = \star \text{ or } x_i = y_i\}$$

We present a construction under the sub-exponential DDH assumption.

**Theorem 2 (Hypercubes from DDH).** *Let  $\mathcal{HC} = \{\mathcal{HC}_{\ell(\lambda)}\}_{\lambda > 0}$  where  $\mathcal{HC}_\ell = \{0, 1, \star\}^\ell$  be the set of hypercubes on  $\{0, 1\}^\ell$ . Then, there is a construction of AGG-PRF supporting the set system  $\mathcal{HC}$  with the product aggregation function, assuming the subexponential DDH assumption.*

We sketch the construction from DDH below. Our DDH construction is the Naor-Reingold PRF [NR04]. Namely, the function is parametrized by an  $\ell$ -tuple  $\mathbf{k} = (k_1, \dots, k_\ell)$  and is defined as

$$F(\mathbf{k}, x) = g^{\prod_{i: x_i=1} k_i}$$

Let us illustrate aggregation over the hypercube  $\mathbf{y} = (1, 0, \star, \star, \dots, \star)$ . To aggregate the function  $F$ , observe that

$$\begin{aligned} \prod_{\{x: x_1=1, x_2=0\}} F(\mathbf{k}, x) &= \prod_{\{x: x_1=1, x_2=0\}} g^{\prod_{i: x_i=1} k_i} \\ &= g^{\sum_{\{x: x_1=1, x_2=0\}} \prod_{i: x_i=1} k_i} \\ &= g^{(k_1)(1)(k_2+1)(k_3+1)\dots(k_\ell+1)} \end{aligned}$$

which can be efficiently computed given  $\mathbf{k}$ .

*Decision Trees.* A decision tree  $T$  on  $\ell$  variables is a binary tree where each internal node is labeled by a variable  $x_i$ , the leaves are labeled by either 0 or 1, one of the two outgoing edges of an internal node is labeled 0, and the other is labeled 1. Computation of a decision tree on an input  $(x_1, \dots, x_\ell)$  starts from the root, and at each internal node  $n$ , proceeds by taking either the 0-outgoing edge or 1-outgoing edge depending on whether  $x_n = 0$  or  $x_n = 1$ , respectively. Finally, the output of the computation is the label of the leaf reached through this process. The size of a decision tree is the number of nodes in the tree.

A decision tree  $T$  defines a set  $S = S_T = \{x \in \{0, 1\}^\ell : T(x) = 1\}$ . We show how to compute product aggregates over sets defined by polynomial size decision trees, under the subexponential DDH assumption.

The construction is simply a result of the observation that the set  $S = S_T$  can be written as a disjoint union of polynomially many hypercubes. Computing aggregates over each hypercube and multiplying the results together gives us the decision tree aggregate.

**Theorem 3 (Decision Trees from DDH).** *Assuming the sub-exponential hardness of the decisional Diffie-Hellman assumption, there is an AGG-PRF that supports aggregation over sets recognized by polynomial-size decision trees.*

*Read-Once Boolean Formulas.* Finally, we show a construction of AGG-PRF over read-once Boolean formulas, the most expressive of our set systems, under the subexponential DDH assumption. A read-once Boolean formula is a Boolean circuit composed of AND, OR and NOT gates with fan-out 1, namely each input literal feeds into at most one gate, and each gate output feeds into at most one other gate. Thus, a read-once formula can be written as a binary tree where each internal node is labeled with an AND or OR gate, and each literal (variable or its negation) appears in at most one leaf.

**Theorem 4 (Read-Once Boolean Formulas from DDH).** *Under the sub-exponential decisional Diffie-Hellman assumption, there is an AGG-PRF that supports aggregation over sets recognized by read-once Boolean formulas.*

Our aggregate PRF is, once again, the Naor-Reingold PRF. The index of the PRF consists of a  $(\ell + 1)$ -tuple of integers in  $\mathbb{Z}_p$ , namely  $\mathbf{K} = (K_0, \dots, K_\ell) \in \mathbb{Z}_p^{\ell+1}$ . The function is defined as

$$f_{\mathbf{K}}(x) = g^{K_0 \prod_{i \in [\ell]} K_i^{x_i}}$$

We compute aggregates by recursion on the levels of the formula. We start by noting that it is enough to compute

$$A(C, 1) := \sum_{x: C(x)=1} \prod_{i \in [1 \dots \ell]} K_i^{x_i}$$

because once this is done, it is easy to compute

$$\prod_{x: C(x)=1} f_{\mathbf{k}}(x) = g^{K_0 \cdot A(C, 1)}$$

For the purposes of this informal exposition, assume that  $\ell$  is a power of two. Let  $C$  be the formula, with either  $C = C_L \wedge C_R$  or  $C = C_L \vee C_R$  for subformula  $C_L$  and  $C_R$ . We show how to recursively compute  $A(C, 1)$  for these sub-circuits and thus for  $C$ .

*Limits of Aggregation.* A natural question to ask is whether one can support aggregation over sets defined by general circuits. It is however easy to see that you cannot support any class of circuits for which deciding satisfiability is hard (for example,  $AC^0$ ), or even ones for which counting the number of SAT assignments is hard (DNFs, for example) as follows. Suppose  $C$  is a circuit which is either unsatisfiable or has a unique SAT assignment. Solving satisfiability for such circuits is known to be sufficient to solve SAT in general [VV86]. The algorithm for SAT simply runs the aggregator with a random PRF key  $K$ , and outputs YES if and only if the aggregator returns a non-zero value. Note that if the formula is unsatisfiable, we will always get 0 from the aggregator. Otherwise, we get  $f_k(x)$ , where  $x$  is the (unique) satisfying assignment. Now, this might end up being 0 accidentally, but cannot be 0 always since otherwise, we will turn it into a PRF distinguisher. The distinguisher has the satisfying assignment hardcoded into it non-uniformly, and it simply checks if  $PRF_K(x)$  is 0.

**Theorem 5 (Impossibility for General Set Systems).** *Suppose there is an efficient algorithm which on an index for  $f \in \mathcal{F}$ , a set system defined by  $\{x : C(x) = 1\}$  for a polynomial size Boolean circuit  $C$ , and an aggregation function  $\Gamma$ , outputs the  $\Gamma_{x:C(x)=1}f(x)$ . Then, there is efficient algorithm that takes circuits  $C$  as input and w.h.p. over its coins, decides satisfiability for  $C$ .*

## 1.2 Our Results: Augmented PRFs and Computational Learning

As discussed above, connections between PRFs and learning theory date back to the 80's in the pioneering work of [Val84] showing that PRF in a complexity class  $C$  implies the existence of concept classes in  $C$  which can not be learned with membership queries. In the second part of this work, we study the implications of the slew of augmented PRF constructions of recent years [BW13a,BGI14a,KPTZ13b,BC10,ABPP14] and our new aggregate PRF to computational learning.

**Constrained PRFs and limits on Restriction Access learnability** Recently, Dvir, Rao, Wigderson, and Yehudayoff [DRWY12] introduced a new learning model where the learner is allowed non-black-box information on the computational device (such as circuits, DNF, formulas) that decides the concept; their learner receives a simplified device resulting from partial assignments to input variables (i.e. restrictions). These partial restrictions lie somewhere in between function evaluation (full restrictions) which correspond to learning with membership queries and the full description of the original device (the empty restriction). The work of [DRWY12] studies a PAC version of restriction access,

called  $\text{PAC}_{RA}$ , where the learner receives the circuit restricted with respect to random partial assignments. They show that both decision trees and DNF formulas can be learned efficiently in this model. Indeed, the  $\text{PAC}_{RA}$  model seems like quite a powerful generalization, if not too unrealistic, of the traditional PAC learning model, as it returns to the learner a computational description of the simplified concept.

Yet, in this section we will show limitations of this computational model under cryptographic assumptions. We show that the *constrained pseudo-random function families* introduced recently in [BW13b,BGI14b,KPTZ13a] naturally define a concept class which is not learnable by an even stronger variant of the restriction access learning model which we define. In the stronger variant, which we name *membership queries with restriction access* ( $\text{MQ}_{RA}$ ) the learner can adaptively specify any restriction of the circuit from a specified class of restrictions  $\mathcal{S}$  and receive the simplified device computing the concept on this restricted domain in return. As this setting requires substantial notation, we define this new model very informally, and defer the formal definitions and theorems to the full version.

**Definition 1 (Membership queries with restriction access ( $\text{MQ}_{RA}$ )).** Let  $\mathcal{C} : X \rightarrow \{0,1\}$  be a concept class, and  $\mathcal{S} = \{S \subseteq X\}$  be a collection of subsets of the domain.  $\mathcal{S}$  is the set of allowable restrictions for concepts  $f \in \mathcal{C}$ . Let  $\text{Simp}$  be “simplification rule” which, for a concept  $f$  and restriction  $S$  outputs a “simplification” of  $f$  restricted to  $S$ .

An algorithm  $\mathcal{A}$  is an  $(\epsilon, \delta, \alpha)$ - $\text{MQ}_{RA}$  learning algorithm for representation class  $\mathcal{C}$  with respect to a restrictions in  $\mathcal{S}$  and simplification rule  $\text{Simp}$  if, for every  $f \in \mathcal{C}$ ,  $\Pr[\mathcal{A}^{\text{Simp}(f,\cdot)} = h] \geq 1 - \delta$  where  $h$  is an  $\epsilon$ -approximation to  $f$  – and furthermore,  $\mathcal{A}$  only requests restrictions for an  $\alpha$ -fraction of the whole domain  $X$ .

Informally, constrained PRFs are PRFs with two additional properties: 1) for any subset  $S$  of the domain in a specified collection  $\mathcal{S}$ , a *constrained key*  $K_S$  can be computed, knowledge of which enables efficient evaluation of the PRF on  $S$ ; and 2) even with knowledge of constrained keys  $K_{S_1}, \dots, K_{S_m}$  for the corresponding subsets, the function retains pseudo-randomness on all points not covered by any of these sets. Connecting this to restriction access, the constrained keys will allow for generation of restriction access examples (restricted implementations with fixed partial assignments) and the second property implies that those examples do not aid in the learning of the function.

**Theorem 6 (Informal).** Suppose  $\mathcal{F}$  is a family of constrained PRFs which can be constrained to sets in  $\mathcal{S}$ . If  $\mathcal{F}$  is computable in circuit complexity class  $\mathcal{C}$ , then  $\mathcal{C}$  is hard to  $\text{MQ}_{RA}$ -learn with restrictions in  $\mathcal{S}$ .

**Corollary 1 (Informal).** Existing constructions of constrained PRFs [BW13a] yield the following corollaries:

- If one-way functions exist, then poly-sized circuits can not be learned with restrictions on sub-intervals of the input-domain; and



- Assuming the sub-exponential hardness of the multi-linear Diffie-Hellman problem,  $NC^1$  cannot be learned with restriction on hypercubes.

**New Learning Models Inspired by the Study of PRFs** We proceed to define two new learning models inspired by recent directions in cryptography. The first model is the *related concept* model inspired by work into related-key attacks in cryptography. While we have cryptography and lower bounds in mind, we argue that this model is in some ways natural. The second model, learning with *aggregate queries*, is directly inspired by our development of aggregate pseudo-random functions in this work; rather than being a natural model in its own right, this model further illustrates how cryptography and learning are duals in many senses.

*The Related Concept Learning Model* The idea that some functions or concepts are related to one another is quite natural. For a DNF formula, for instance, related concepts may include formulas where a clause has been added or formulas where the roles of two variables are swapped. For a decision tree, we could consider removing some accepting leaves and examining the resulting behavior. For a circuit, a related circuit might alter internal gates or fix the values on some wires. A similar phenomena occurs in cryptography, where secret keys corresponding to different instances of the same cryptographic primitive or even secret keys of different cryptographic primitives are related (if, for example, they were generated by a pseudo random process on the same seed).

We propose a new computational learning model where the learner is explicitly allowed to specify membership queries not only for the concept to be learned, but also for “related” concepts, given by a class of allowed transformations on the concept. We will show both a separation from membership queries, and a general negative result in the new model. Based on recent constructions of related-key secure PRFs by Bellare and Cash [BC10] and Abdalla et al [ABPP14], we demonstrate concept classes for which access to these related concepts is of no help.

To formalize the related concept learning model, we will consider *keyed concept classes* – classes indexed by a set of keys. This will enable the study of related concepts by instead considering concepts whose keys are related in some way. Most generally, we think of a key as a succinct representation of the computational device which decides the concept. This is a general framework; for example, we may consider the bit representation of a particular log-depth circuit as a key for a concept in the concept class  $NC^1$ . For a concept  $f_k$  in concept class  $\mathcal{C}$ , we allow the learner to query a membership oracle for  $f_k$  and also for ‘related’ concepts  $f_{\phi(k)} \in \mathcal{C}_K$  for  $\phi$  in a specified class of allowable functions  $\Phi$ . For example: let  $K = \{0, 1\}^\lambda$  and let  $\Phi^\oplus = \{\phi_\Delta : k \mapsto k \oplus \Delta\}_{\Delta \in \{0, 1\}^\lambda}$ . Informally:

**Definition 2 ( $\Phi$ -Related-Concept Learning Model ( $\Phi$ -RC)).** For  $\mathcal{C}_K$  a keyed concept class, let  $\Phi = \{\phi : K \rightarrow K\}$  be a set of functions on  $K$  that contains the identity function  $\text{id}$ . A related-concept oracle  $RC_k$ , on query  $(\phi, x)$ , responds with  $f_{\phi(k)}(x)$ , for all  $\phi \in \Phi$  and  $x \in X$ .

An algorithm  $A$  is an  $(\epsilon, \delta)$ - $\Phi$ -RK learning algorithm for a  $\mathcal{C}_k$  if, for every  $k \in K$ , when given access to the oracle  $RK_k(\cdot)$ , the algorithm  $A$  outputs with probability at least  $1 - \delta$  a function  $h : \{0, 1\}^n \rightarrow \{0, 1\}$  that  $\epsilon$ -approximates  $f_k$ .

Yet again, we are able to demonstrate the limitations of this model using the power of a strong type of pseudo-random function. We show that *related-key secure PRF families (RKA-PRF)* defined and instantiated in [BC10] and [ABPP14] give a natural concept class which is not learnable with related key queries. RKA-PRFs are defined with respect to a set  $\Phi$  of functions on the set of PRF keys. Informally, the security notion guarantees that for a randomly selected key  $k$ , no efficient adversary can distinguish oracle access to  $f_k$  and  $f_{\phi(k)}$  (for many adaptively chosen functions  $\phi \in \Phi$ ) from an oracle that returns completely random values. We leverage this strong pseudo-randomness property to show hard-to-learn concepts in the related concept model.

**Theorem 7 (Informal).** *Suppose  $\mathcal{F}$  is a family of RKA-PRFs with respect to related-key functions  $\Phi$ . If  $\mathcal{F}$  is computable in circuit complexity class  $\mathcal{C}$ , then  $\mathcal{C}$  is hard to learn in the  $\Phi'$ -RC model for some  $\Phi'$ .*

Existing constructions of RKA-PRFs [ABPP14] yield the following corollary:

**Corollary 2 (Informal).** *Assuming the hardness of the DDH problem, and collision-resistant hash functions,  $NC^1$  is hard to  $\Phi$ -RC-learn for an class of affine functions  $\Phi$ .*

*The Aggregate Learning Model* The other learning model we propose is inspired by our aggregate PRFs. Here, we consider a new extension to the power of the learning algorithm. Whereas membership queries are of the form “What is the label of an example  $x$ ?”, we grant the learner the power to request the evaluation of simple functions on tuples of examples  $(x_1, \dots, x_n)$  such as “How many of  $x_1, \dots, x_n$  are in  $C$ ?” or “Compute the product of the labels of  $x_1, \dots, x_n$ ?”. Clearly, if  $n$  is polynomial then this will result only a polynomial gain in the query complexity of a learning algorithm in the best case. Instead, we propose to study cases when  $n$  may be super-polynomial, but the description of the tuples is succinct. For example, the learning algorithm might query the number of  $x$ ’s in a large interval that are positive examples in the concept.

As with the restriction access and related concept models – and the aggregate PRFs we define in this work – the Aggregate Queries (AQ) learning model will be considered with restrictions to both the types of aggregate functions  $\Gamma$  the learner can query, and the sets  $\mathcal{S}$  over which the learner may request these functions to be evaluated on. We now present the AQ learning model informally:

**Definition 3 (( $\Gamma, \mathcal{S}$ )-Aggregate Queries (AQ) Learning).** *Let  $\mathcal{C} : X \rightarrow \{0, 1\}$  be a concept class, and let  $\mathcal{S}$  be a collection of subsets of  $X$ . Let  $\Gamma : \{0, 1\}^* \rightarrow V$  be an aggregation function. For  $f \in \mathcal{C}$ , let  $\text{AGG}_f$  be an “aggregation” oracle, which for  $S \in \mathcal{S}$ , returns  $\Gamma_{x \in S} f(x)$ . Let  $\text{MEM}_f$  be the membership oracle, which for input  $x$  returns  $f(x)$ .*

An algorithm  $\mathcal{A}$  is an  $(\epsilon, \delta)$ - $(\Gamma, \mathcal{S})$ -AQ learning algorithm for  $\mathcal{C}$  if for every  $f \in \mathcal{C}$ ,

$$\Pr[\mathcal{A}^{MEM_f(\cdot), AGG_f(\cdot)} = h] \geq 1 - \delta$$

where  $h$  is an  $\epsilon$ -approximation to  $f$ .

Initially, AQ learning is reminiscent of learning with statistical queries (SQ). In fact, this apparent connection inspired this portion of our work. But the AQ setting is in fact incomparable to SQ learning, or even the weaker “statistical queries that are independent of the target” as defined in [BF02]. On the one hand, AQ queries provide a sort of noiseless variant of SQ, giving more power to the AQ learner; on the other hand, the AQ learner is restricted to aggregating over sets in  $\mathcal{S}$ , whereas the SQ learner is not restricted in this way, thereby limiting the power of the AQ learner. The AQ setting where  $\mathcal{S}$  contains every subset of the domain is indeed a noiseless version of “statistical queries independent of the target,” but even this model is a restricted version of SQ. This does raise the natural question of a noiseless version of SQ and its variants; hardness results in such models would be interesting in that they would suggest that the hardness comes not from the noise but from an inherent loss of information in statistics/aggregates.

We will show both a simple separation from learning with membership queries (in the full version), and under cryptographic assumptions, a general lower bound on the power of learning with aggregate queries. The negative examples will use the results in Section 1.1.

**Theorem 8.** *Let  $\mathcal{F}$  be a boolean-valued aggregate PRF with respect to set system  $\mathcal{S}$  and aggregation function  $\Gamma$ . If  $\mathcal{F}$  is computable in complexity class  $\mathcal{C}$ , then  $\mathcal{C}$  is hard to  $(\Gamma, \mathcal{S})$ -AQ learn.*

**Corollary 3.** *Using the results from Section 3, we get the following corollaries:*

- *The existence of one way functions implies that  $P/poly$  is hard to  $(\sum, \mathcal{S}_{[a,b]})$ -AQ learn, with  $\mathcal{S}_{[a,b]}$  the set of sub-intervals of the domain as defined in Section 3.*
- *The DDH assumption implies that  $NC^1$  is hard to  $(\sum, \mathcal{S}_{[a,b]})$ -AQ learn, with  $\mathcal{S}_{[a,b]}$  being the set of sub-intervals of the domain as defined in Section 3.*
- *The subexponential DDH Assumption implies that  $NC^1$  is hard to  $(\prod, \mathcal{R})$ -AQ learn, with  $\mathcal{R}$  the set of read-once boolean formulas defined in Section 3.*

*Open Questions.* As discussed in the introduction, augmented pseudo-random functions often have powerful and surprising applications, perhaps the most recent example being constrained PRFs [BW13a, KPTZ13a, BGI14a]. Perhaps the most obvious open question that emerges from this work is to find applications for aggregate PRFs. We remark that a primitive similar to aggregate PRFs was used in [BGV11] to construct delegation protocols.

Perhaps a more immediate concern is that all our aggregate PRF constructions (except for intervals) requires sub-exponential hardness assumptions. We

view it as an important open question to base these constructions on polynomial assumptions.

In this work we restricted our attention to particular types of aggregation functions and subsets over which the aggregation takes place, although our definition captures more general scenarios. We looked at aggregation functions that compute group operations over Abelian groups. Can we support more general aggregation functions that are not restricted to group operations, for example the majority aggregation function, or even non-symmetric aggregation functions? We show positive results for intervals, hypercubes, and sets recognized by read-once formulas and decision trees. On the other hand, we show that it is unlikely that we can support general sets, for example sets recognized by CNF formulas. This almost closes the gap between what is possible and what is hard. A concrete open question in this direction is to construct an aggregate PRF computing summation over an Abelian group for sets recognized by DNFs, or provide evidence that this cannot be done.

*Organization.* This paper is organized into two parts that can be read essentially independently of each other. In the first part (Sections 2 and 3), we present the definition and constructions of aggregate pseudo-random functions. In the second part (Section 4), we show connections between various notions of augmented PRFs and their applications to augmented learning models.

## 2 Aggregate PRF

We will let  $\lambda$  denote the security parameter throughout this paper.

Let  $\mathcal{F} = \{\mathcal{F}_\lambda\}_{\lambda>0}$  be a function family where each function  $f \in \mathcal{F}_\lambda$  maps a domain  $\mathcal{D}_\lambda$  to a range  $\mathcal{R}_\lambda$ . An *aggregate function* family is associated with two objects:

1. an ensemble of sets  $\mathcal{S} = \{\mathcal{S}_\lambda\}_{\lambda>0}$  where each  $\mathcal{S}_\lambda$  is a collection of subsets of the domain  $S \subseteq \mathcal{D}_\lambda$ ; and
2. an “aggregation function”  $\Gamma_\lambda : (\mathcal{R}_\lambda)^* \rightarrow \mathcal{V}_\lambda$  that takes a tuple of values from the range  $\mathcal{R}_\lambda$  of the function family and “aggregates” them to produce a value in an output set  $\mathcal{V}_\lambda$ .

Let us now make this notion formal. To do so, we will impose restrictions on the set ensembles and the aggregation function. First, we require set ensemble  $\mathcal{S}_\lambda$  to be *efficiently recognizable*. That is, there is a polynomial-size Boolean circuit family  $\mathcal{C} = \{\mathcal{C}_\lambda\}_{\lambda>0}$  such that for any set  $S \in \mathcal{S}_\lambda$  there is a circuit  $C = C_S \in \mathcal{C}_\lambda$  such that  $x \in S$  if and only if  $C(x) = 1$ . Second, we require our aggregation functions  $\Gamma$  to be efficient in the length of its inputs, and symmetric; namely the output of the function does not depend on the order in which the inputs are fed into it. Summation over an Abelian group is an example of a possible aggregation function. Third and finally, elements in our sets  $\mathcal{D}_\lambda$ ,  $\mathcal{R}_\lambda$ , and  $\mathcal{V}_\lambda$  are all representable in  $\text{poly}(\lambda)$  bits, and the functions  $f \in \mathcal{F}_\lambda$  are computable in  $\text{poly}(\lambda)$  time.

Define the aggregate function  $\text{AGG} = \text{AGG}_{f, \mathcal{S}_\lambda, \Gamma_\lambda}^\lambda$  that is indexed by a function  $f \in \mathcal{F}_\lambda$ , takes as input a set  $S \in \mathcal{S}_\lambda$  and “aggregates” the values of  $f(x)$  for all  $x \in \mathcal{S}_\lambda$ . That is,  $\text{AGG}(S)$  outputs

$$\Gamma(f(x_1), f(x_2), \dots, f(x_{|S|}))$$

where  $S = \{x_1, \dots, x_{|S|}\}$ . More precisely, we have

$$\begin{aligned} \text{AGG}_{f, \mathcal{S}_\lambda, \Gamma_\lambda}^\lambda : \mathcal{S}_\lambda &\rightarrow \mathcal{V}_\lambda \\ S &\mapsto \Gamma_{x_i \in S}(f(x_1), \dots, f(x_{|S|})) \end{aligned}$$

We will furthermore require that the AGG can be computed in  $\text{poly}(\lambda)$  time. We require this in spite of the fact that the sets over which the aggregation is done can be exponentially large! Clearly, such a thing is impossible for a random function  $f$  but yet, we will show how to construct *pseudo-random* function families that support efficient aggregate evaluation. We will call such a pseudo-random function (PRF) family an *aggregate PRF* family. In other words, our objective is two fold:

1. Allow anyone who knows the (polynomial size) function description to efficiently compute the aggregate function values over exponentially large sets; but at the same time,
2. Ensure that the function family is indistinguishable from a truly random function, even given an oracle that computes aggregate values.

A simple example of aggregates is that of computing the summation of function values over sub-intervals of the domain. That is, let domain and range be  $\mathbb{Z}_p$  for some  $p = p(\lambda)$ , let the family of subsets be  $\mathcal{S}_\lambda = \{[a, b] \subseteq \mathbb{Z}_p : a, b \in \mathbb{Z}_p; a \leq b\}$ , and the aggregation function be  $\Gamma_\lambda(y_1, \dots, y_k) = \sum_{i=1}^k y_i \pmod{p}$ . In this case, we are interested in computing

$$\text{AGG}_{f, \mathcal{S}_\lambda, \text{sum}}^\lambda([a, b]) = \sum_{a \leq x \leq b} f(x)$$

We will, in due course, show both constructions and impossibility results for aggregate PRFs, but first let us start with the formal definition.

**Definition 4 (Aggregate PRF).** *Let  $\mathcal{F} = \{\mathcal{F}_\lambda\}_{\lambda > 0}$  be a function family where each function  $f \in \mathcal{F}_\lambda$  maps a domain  $\mathcal{D}_\lambda$  to a range  $\mathcal{R}_\lambda$ ,  $\mathcal{S}$  be an efficiently recognizable ensemble of sets  $\{\mathcal{S}_\lambda\}_{\lambda > 0}$ , and  $\Gamma_\lambda : (\mathcal{R}_\lambda)^* \rightarrow \mathcal{V}_\lambda$  be an aggregation function. We say that  $\mathcal{F}$  is an  $(\mathcal{S}, \Gamma)$ -aggregate pseudorandom function family (also denoted  $(\mathcal{S}, \Gamma)$ -AGG-PRF) if there exists an efficient algorithm  $\text{Aggregate}_{k, \mathcal{S}, \Gamma}(S)$ : On input a subset  $S \in \mathcal{S}$  of the domain, outputs  $v \in \mathcal{V}$ , such that*

- **Efficient aggregation:** For every  $S \in \mathcal{S}$ ,  $\text{Aggregate}_{k, \mathcal{S}, \Gamma}(S) = \text{AGG}_{k, \mathcal{S}, \Gamma}(S)$  where  $\text{AGG}_{k, \mathcal{S}, \Gamma}(S) := \Gamma_{x \in S} F_k(x)$ .<sup>67</sup>

<sup>6</sup> We omit subscripts on AGG and Aggregate when clear from context.

<sup>7</sup> AGG is defined to be the correct aggregate value, while Aggregate is the algorithm by which we compute the value AGG. We make this distinction because while a random function cannot be efficiently aggregated, the aggregate value is still well-defined.

- **Pseudorandomness:** For all probabilistic polynomial-time (in security parameter  $\lambda$ ) algorithms  $A$ , and for randomly selected key  $k \in K$ :

$$\left| \Pr_{f \leftarrow \mathcal{F}_\lambda} [A^{f_k, AGG_{f_k, s, r}}(1^\lambda)] - \Pr_{h \leftarrow \mathcal{H}_\lambda} [A^{h, AGG_{h, s, r}}(1^\lambda)] \right| \leq \text{negl}(\lambda)$$

where  $\mathcal{H}_\lambda$  is the set of all functions  $D_\lambda \rightarrow R_\lambda$ .

*Remark 1.* In this work, we restrict our attention to aggregation functions that treat the range  $\mathcal{V}_\lambda = \mathcal{R}_\lambda$  as an Abelian group and compute the group sum (or product) of its inputs. We denote this setting by  $\Gamma = \sum$  (or  $\prod$ , respectively). Supporting other types of aggregation functions (ex: max, a hash) is a direction for future work.

## 2.1 A General Security Theorem for Aggregate PRFs

How does the security of a function family in the AGG-PRF game relate to security in the normal PRF game (in which  $A$  uses only the oracle  $f$  and not  $AGG_f$ )?

In this section, we show a general security theorem for aggregate pseudo-random functions. Namely, we show that any “sufficiently secure” PRF is also aggregation-secure (for any collection of efficiently recognizable sets and any group-aggregation operation), in the sense of Definition 4, by way of an *inefficient* reduction (with overhead polynomial in the size of the domain). In Section 3, we will use this to construct AGG-PRFs from a subexponential-time hardness assumption on the DDH problem. We also show that no such *general* reduction can be efficient, by demonstrating a PRF family that is not aggregation-secure. As a general security theorem cannot be shown without the use of complexity leveraging, this suggests a natural direction for future study: to devise constructions for similarly expressive aggregate PRFs from polynomial assumptions.

**Lemma 1.** *Let  $\mathcal{F} = \{\mathcal{F}_\lambda\}_{\lambda > 0}$  be a pseudo-random function family where each function  $f \in \mathcal{F}_\lambda$  maps a domain  $\mathcal{D}_\lambda$  to a range  $\mathcal{R}_\lambda$ . Suppose there is an adversary  $A$  that runs in time  $t_A = t_A(\lambda)$  and achieves an advantage of  $\epsilon_A = \epsilon_A(\lambda)$  in the aggregate PRF security game for the family  $\mathcal{F}$  with an efficiently recognizable set system  $\mathcal{S}_\lambda$  and an aggregation function  $\Gamma_\lambda$  that is computable in time polynomial in its input length. Then, there is an adversary  $B$  that runs in time  $t_B = t_A + \text{poly}(\lambda, |\mathcal{D}_\lambda|)$  and achieves an advantage of  $\epsilon_B = \epsilon_A$  in the standard PRF game for the family  $\mathcal{F}$ .*

*Proof.* Let  $f_K \leftarrow \mathcal{F}_\lambda$  be a random function from the family  $\mathcal{F}_\lambda$ . We construct the adversary  $B$  which is given access to an oracle  $\mathcal{O}$  which is either  $f_K$  or a uniformly random function  $h : \mathcal{D}_\lambda \rightarrow \mathcal{R}_\lambda$ .

$B$  works as follows: It queries the PRF on all inputs  $x \in \mathcal{D}_\lambda$ , builds the function table  $T_K$  of  $f_K$  and runs the adversary  $A$ , responding to its queries as follows:

1. Respond to its PRF query  $x \in \mathcal{D}_\lambda$  by returning  $T_K[x]$ ; and

2. Respond to its aggregate query  $(I, S)$  by (a) going through the table to look up all  $x$  such that  $x \in S$ ; and (b) applying the aggregation function honestly to these values.

Finally, when  $A$  halts and returns a bit  $b$ ,  $B$  outputs the bit  $b$  and halts.  $B$  takes  $O(|\mathcal{D}_\lambda|)$  time to build the truth table of the oracle. For each aggregate query  $(I, S)$ ,  $B$  first checks for each  $x \in \mathcal{D}_\lambda$  whether  $x \in S$ . This takes  $|\mathcal{D}_\lambda| \cdot \text{poly}(\lambda)$  time, since  $S$  is efficiently recognizable. It then computes the aggregation function  $I$  over  $f(x)$  such that  $x \in S$ , taking  $\text{poly}(|\mathcal{D}_\lambda|)$  time, since  $I$  is computable in time polynomial in its input length. The total time, therefore, is

$$t_B = t_A + \text{poly}(\lambda, |\mathcal{D}_\lambda|)$$

Clearly, when  $\mathcal{O}$  is the pseudo-random function  $f_K$ ,  $B$  simulates an aggregatable PRF oracle to  $A$ , and when  $\mathcal{O}$  is a random function,  $B$  simulates an aggregate random oracle to  $A$ . Thus,  $B$  has the same advantage in the PRF game as  $A$  does in the aggregate PRF game.

The above gives an inefficient reduction from the PRF security of a function family  $\mathcal{F}$  to the AGG-PRF security of the same family running in time polynomial in the size of the domain. Can this reduction be made efficient; that is, can we replace  $t_B = t_A + \text{poly}(\lambda)$  into the Lemma 1?

This is not possible. Such a reduction would imply that every PRF family that supports efficient aggregate functionality AGG is AGG-PRF secure; this is clearly false. Take for example a pseudorandom function family  $\mathcal{F}_0 = \{f : \mathbb{Z}_{2p} \rightarrow \mathbb{Z}_p\}$  such that for all  $f$ , there is no  $x$  with  $f(x) = 0$ . It is possible to construct such a pseudorandom function family  $\mathcal{F}_0$  (under the standard definition). While 0 is not in the image of any  $f \in \mathcal{F}_0$ , a random function with the same domain and range will, with high probability, have 0 in the image. For an aggregation oracle  $\text{AGG}_f$  computing *products* over  $\mathbb{Z}_p$ :  $\text{AGG}_f(\mathbb{Z}_{2p}) \neq 0$  if  $f \in \mathcal{F}_0$ , while  $\text{AGG}_f(\mathbb{Z}_{2p}) = 0$  with high probability for random  $f$ .

Thus, access to aggregates for products over  $\mathbb{Z}_p$ <sup>8</sup> would allow an adversary to trivially distinguish  $f \in \mathcal{F}_0$  from a truly random map.

## 2.2 Impossibility of Aggregate PRF for General Sets

It is natural to ask whether whether an aggregate PRF might be constructed for more general sets than we present in Section 3. There we constructed aggregate PRF for the sets of all satisfying assignments for read-once boolean formula and decision trees. As we show in the following, it is impossible to extend this to support the set of satisfying assignments for more general circuits.

<sup>8</sup> Taken with respect to a set ensemble  $\mathcal{S}$  containing, as an element, the whole domain  $\mathbb{Z}_{2p}$ . While this is not necessary (a sufficiently large subset would suffice), it is the case for the ensembles  $\mathcal{S}$  we consider in this work.

**Theorem 9.** *Suppose there is an algorithm that has a PRF description  $K$ , a circuit  $C$ , and a fixed aggregation rule (sum over a finite field, say), and outputs the aggregate value*

$$\sum_{x:C(x)=1} f_K(x)$$

*Then, there is an algorithm that takes circuits  $C$  as input and w.h.p. over its coins, decides the satisfiability of  $C$ .*

*Proof.* The algorithm for SAT simply runs the aggregator with a randomly chosen  $K$ , and outputs YES if and only if the aggregator returns 1. The rationale is that if the formula is unsatisfiable, you will always get 0 from the aggregator.<sup>9</sup> Otherwise, you will get  $f_K(x)$ , where  $x$  is the satisfying assignment. (More generally,  $\sum_{x:C(x)=1} f_K(x)$ ). Now, this might end up being 0 accidentally, but cannot be 0 always since otherwise, you will get a PRF distinguisher. The distinguisher has the satisfying assignment hardcoded into it non-uniformly,<sup>10</sup> and it simply checks if  $f_K(x) = 0$ .

This impossibility result can be generalized for efficient aggregation of functions that are not pseudo-random. For instance, if  $f(x) \equiv 1$  was the constant function 1, the same computing the aggregate over  $f$  satisfying inputs to  $C$  would not only reveal the satisfiability of  $C$ , but even the number of satisfying assignments! In the PRF setting though, it seems that aggregates only reveal the (un)satisfiability of a circuit  $C$ , but not the number of satisfying assignments. Further studying the relationship between the (not necessarily pseudo-random) function  $f$ , the circuit representation of  $C$ , and the tractability of computing aggregates is an interesting direction. A negative result for a class for which satisfiability (or even counting assignments) is tractable would be very interesting.

### 3 Constructions of aggregate PRF

In this section, we show several constructions of aggregate PRFs. In Section 3.1, we show a generic construction of aggregate PRFs for intervals (where the aggregation is any group operation). This construction is black-box: given any PRF with the appropriate domain and range, we construct a related family of aggregate PRFs and with no loss in security. In Section 3.2, we show a construction of aggregate PRFs for products over bit-fixing sets (hypercubes), from a strong decisional Diffie-Hellman assumption. We then generalize the DDH construction:

<sup>9</sup> This proof may be extended to the case when the algorithm's output is not restricted to be 0 when the input circuit  $C$  is unsatisfiable, and even arbitrary outputs for sufficiently expressive classes of circuits.

<sup>10</sup> As pointed out by one reviewer, for sufficiently expressive classes of circuits  $C$ , this argument can be made uniform. Specifically, we use distinguish the challenge  $y$  from a pseudo-random generator from random by choosing  $C := C_y$  that is satisfiable if and only if  $y$  is in the PRG image, and modify the remainder of the argument accordingly.



in Section 3.3, to the class of sets recognized by polynomial-size decision trees; and in Section 3.4, to sets recognized by read-once Boolean formulas. In these last three constructions, we make use of Lemma 1 to argue security.

### 3.1 Generic Construction for Interval Sets

Our first construction is adapted from [GGN10]<sup>11</sup>. The construction is entirely black-box: from any appropriate PRF family  $\mathcal{G}$ , we construct a related AGG-PRF family  $\mathcal{F}$ . Unlike the proofs in the sequel, this reduction exactly preserves the security of the starting PRF.

Let  $\mathcal{G}_\lambda = \{g_K : \mathbb{Z}_{n(\lambda)} \rightarrow R_\lambda\}_{K \in \mathcal{K}_\lambda}$  be a PRF family, with  $R = R_\lambda$  being a group where the group operation is denoted by  $\oplus$ <sup>12</sup>. We construct an aggregatable PRF  $\mathcal{F}_\lambda = \{f_K\}_{K \in \mathcal{K}_\lambda}$  for which we can efficiently compute summation of  $f_K(x)$  for all  $x$  in an interval  $[a, b]$ , for any  $a \leq b \in \mathbb{Z}_n$ . Let  $\mathcal{S}_{[a,b]} = \{[a, b] \subseteq \mathbb{Z}_n : a, b \in \mathbb{Z}_n; a \leq b\}$  be the set of all interval subsets of  $\mathbb{Z}_n$ ,  $[a, b] = \{x \in \mathbb{Z}_n : a \leq x \leq b\}$ . Define  $\mathcal{F} = \{f_K : \mathbb{Z}_n \rightarrow R\}_{K \in \mathcal{K}}$  as follows:

$$f_K(x) = \begin{cases} g_K(0) & : x = 0 \\ g_K(x) \oplus g_K(x-1) & : x \neq 0 \end{cases}$$

**Lemma 2.** *Assuming that  $\mathcal{G}$  is a pseudo-random function family,  $\mathcal{F}$  is a  $(\mathcal{S}_{[a,b]}, \oplus)$ -aggregate pseudo-random function family.*

*Proof.* It follows immediately from the definition of  $f_K$  that one can compute the summation of  $f_K(x)$  over any interval  $[a, b]$ . Indeed, rearranging the definition yields

$$\sum_{x \in [0, b]} f_K(x) = g_K(b) \quad \text{and} \quad \sum_{x \in [a, b]} f_K(x) = g_K(b) \oplus -g_K(a-1)$$

We reduce the pseudo-randomness of  $\mathcal{F}$  to that of  $\mathcal{G}$ . The key observation is that each query to the  $f_K$  oracle as well as the aggregation oracle for  $f_K$  can be answered using at most two black-box calls to the underlying function  $g_K$ . By assumption on  $\mathcal{G}$ , replacing the oracle for  $g_K$  with a uniformly random function  $h : \mathbb{Z}_n \rightarrow R$  is computationally indistinguishable. Furthermore, the function  $f$  defined by replacing  $g$  by  $h$ , namely

$$f'(x) = \begin{cases} h(0) & : x = 0 \\ h(x) \oplus h(x-1) & : x \neq 0 \end{cases}$$

is a truly random function. Thus, the simulated oracle with  $g_K$  replaced by  $h$  implements a uniformly random function that supports aggregate queries. Security according to Definition 4 follows immediately.

<sup>11</sup> See Example 3.1 and Footnote 18

<sup>12</sup> The only structure of  $\mathbb{Z}_n$  we use is the *total order*. Our construction directly applies to any finite, totally-ordered domain  $D$  by first mapping  $D$  to  $\mathbb{Z}_n$ , preserving order.

### 3.2 Bit-Fixing Aggregate PRF from DDH

We now construct an aggregate PRF computing products for bit-fixing sets. Informally, our PRF will have domain  $\{0, 1\}^{\text{poly}(\lambda)}$ , and support aggregation over sets like  $\{x : x_1 = 0 \wedge x_2 = 1 \wedge x_7 = 0\}$ . We will naturally represent such sets by a string in  $\{0, 1, \star\}^{\text{poly}(\lambda)}$  with 0 and 1 indicating a fixed bit location, and  $\star$  indicating a free bit location. We call each such set a ‘hypercube.’ The PRF will have a multiplicative group  $\mathcal{G}$  as its range, and the aggregate functionality will compute group products.

Our PRF is exactly the Naor-Reingold PRF [NR04], for which we demonstrate efficient aggregation and security. We begin by stating the decisional Diffie-Hellman assumption.

Let  $\mathcal{G} = \{\mathcal{G}_\lambda\}_{\lambda>0}$  be a family of groups of order  $p = p(\lambda)$ . The decisional Diffie-Hellman assumption for  $\mathcal{G}$  says that the following two ensembles are computationally indistinguishable:

$$\begin{aligned} & \{(\mathcal{G}_\lambda, g, g^a, g^b, g^{ab}) : G \leftarrow \mathcal{G}_\lambda; g \leftarrow G; a, b \leftarrow \mathbb{Z}_p\}_{\lambda>0} \\ & \approx_c \{(G, g, g^a, g^b, g^c) : G \leftarrow \mathcal{G}_\lambda; g \leftarrow G; a, b, c \leftarrow \mathbb{Z}_p\}_{\lambda>0} \end{aligned}$$

We say that the  $(t(\lambda), \epsilon(\lambda))$ -DDH assumption holds if for every adversary running in time  $t(\lambda)$ , the advantage in distinguishing between the two distributions above is at most  $\epsilon(\lambda)$ .

**Construction** Let  $\mathcal{G} = \{\mathcal{G}_\lambda\}_{\lambda>0}$  be a family of groups of order  $p = p(\lambda)$ , each with a canonical generator  $g$ , for which the decisional Diffie Hellman (DDH) problem is hard. Let  $\ell = \ell(\lambda)$  be a polynomial function. We will construct a PRF family  $\mathcal{F}_\ell = \{\mathcal{F}_{\ell, \lambda}\}_{\lambda>0}$  where each function  $f \in \mathcal{F}_{\ell, \lambda}$  maps  $\{0, 1\}^{\ell(\lambda)}$  to  $\mathcal{G}_\lambda$ . Our PRF family is exactly the Naor-Reingold PRF [NR04]. Namely, each function  $f$  is parametrized by  $\ell + 1$  numbers  $\mathbf{K} := (K_0, K_1, \dots, K_\ell)$ , where each  $K_i \in \mathbb{Z}_p$ .

$$f_{\mathbf{K}}(x_1, \dots, x_\ell) = g^{K_0 \prod_{i=1}^{\ell} K_i^{x_i}} = g^{K_0 \prod_{i: x_i=1} K_i} \in \mathcal{G}_\lambda$$

The aggregation algorithm **Aggregate** for bit-fixing functions gets as input the PRF key  $\mathbf{K}$  and a bit-fixing string  $y \in \{0, 1, \star\}^\ell$  and does the following:

- Define the strings  $K'_i$  as follows:

$$K'_i = \begin{cases} 1 & \text{if } y_i = 0 \\ K_i & \text{if } y_i = 1 \\ 1 + K_i & \text{otherwise} \end{cases}$$

- Output  $g^{K_0 \prod_{i=1}^{\ell} K'_i}$  as the answer to the aggregate query.

Letting  $\mathcal{HC} = \{\mathcal{HC}_{\ell(\lambda)}\}_{\lambda>0}$  where  $\mathcal{HC}_\ell = \{0, 1, \star\}^\ell$  is the set of hypercubes on  $\{0, 1\}^\ell$ , we now prove the following:

**Theorem 10.** *Let  $\epsilon > 0$  be a constant, choose the security parameter  $\lambda = \Omega(\ell^{1/\epsilon})$ , and assume the  $(2^{\lambda^\epsilon}, 2^{-\lambda^\epsilon})$ -hardness of DDH over the group  $\mathcal{G}$ . Then, the collection of functions  $\mathcal{F}$  defined above is a secure aggregate PRF with respect to the subsets  $\mathcal{HC}$  and the product aggregation function over  $\mathcal{G}$ .*

*Correctness.* We show that the answer we computed for an aggregate query  $y \in \{0, 1, \star\}^\lambda$  is correct. Define the sets

$$\text{Match}(y) := \{x \in \{0, 1\}^\lambda : \forall i, y_i = \star \text{ or } x_i = y_i\} \text{ and } \text{Fixed}(y) := \{i \in [\lambda] : y_i \in \{0, 1\}\}$$

Thus,  $\text{Match}(y)$  is the set of all 0-1 strings  $x$  that match all the fixed locations of  $y$ , but can take any value on the wildcard locations of  $y$ .  $\text{Fixed}(y)$  is the set of all locations  $i$  where the bit  $y_i$  is fixed. Note that:

$$\begin{aligned} \text{AGG}(\mathbf{K}, y) &= \prod_{x \in \text{Match}(y)} f_{\mathbf{K}}(x) && \text{(by definition of AGG)} \\ &= \prod_{x \in \text{Match}(y)} g^{K_0 \prod_{i=1}^\ell K_i^{x_i}} && \text{(by definition of } f_{\mathbf{K}}\text{)} \\ &= g^{K_0 \sum_{x \in \text{Match}(y)} \prod_{i=1}^\ell K_i^{x_i}} \\ &= g^{K_0 \left( \prod_{i \in \text{Fixed}(y)} K_i^{y_i} \right) \cdot \left( \prod_{i \in [\ell] \setminus \text{Fixed}(y)} (1 + K_i) \right)} && \text{(inverting sums and products)} \\ &= g^{K_0 \prod_{i=1}^\ell K_i'} && \text{(by definition of } K_i'\text{)} \\ &= \text{Aggregate}(\mathbf{K}, y) && \text{(by definition of Aggregate)} \end{aligned}$$

*Security.* We will rely on the following theorem from [NR04].

**Theorem 11 (Theorem 4.1, [NR04]).** *Suppose there is an adversary  $A$  that runs in time  $t(\lambda)$  and has an advantage of  $\gamma(\lambda)$  in the (regular) PRF game. Then, there is an adversary  $B$  that runs in time  $\text{poly}(\lambda) \cdot t(\lambda)$  and breaks the DDH assumption with advantage  $\gamma(\lambda)/\lambda$ .*

The aggregate PRF security proof proceeds as follows. First, we choose the security parameter  $\lambda = \Omega(\ell^{1/\epsilon})$  as in the theorem statement. We use Lemma 1 to conclude that if there is an adversary distinguisher  $D$  breaking the aggregate PRF security of  $\mathcal{F}$  in  $\text{poly}(\lambda)$  time with  $1/\text{poly}(\lambda)$  advantage, then there is an adversary  $A$  that breaks the regular PRF security of  $\mathcal{F}$  in  $\text{poly}(\lambda) \cdot 2^{O(\ell)} = \text{poly}(\lambda) \cdot 2^{\lambda^\epsilon} = 2^{O(\lambda^\epsilon)}$  time with  $1/\text{poly}(\lambda)$  advantage. Using Theorem 11 now tells us that there is an adversary  $B$  that wins the DDH distinguishing game in  $2^{O(\lambda^\epsilon)}$  time with  $1/\text{poly}(\lambda)$  advantage, breaking the subexponential DDH assumption. This establishes the aggregate security of the PRF and thus Theorem 10.

Obtaining a security proof based on polynomial assumptions is an interesting open question.

### 3.3 Decision Trees

We generalize the previous construction from DDH to support sets specified by polynomial-sized decision trees by observing that such decision trees can be written as disjoint unions of hypercubes.

A decision tree family  $\mathcal{T}_\lambda$  of size  $p(\lambda)$  over  $\ell(\lambda)$  variables consists of binary trees with at most  $p(\lambda)$  nodes, where each internal node is labeled with a variable

$x_i$  for  $i \in [\ell]$ , the two outgoing edges of an internal node are labeled 0 and 1, and the leaves are labeled with 0 or 1. On input an  $x \in \{0, 1\}^\ell$ , the computation of the decision tree starts from the root, and upon reaching an internal node  $n$  labeled by a variable  $x_i$ , takes either the 0-outgoing edge or the 1-outgoing edge out of the node  $n$ , depending on whether  $x_i$  is 0 or 1, respectively.

We now show how to construct a PRF family  $\mathcal{F}_\ell = \{\mathcal{F}_{\ell, \lambda}\}_{\lambda > 0}$  where each  $\mathcal{F}_{\ell, \lambda}$  consists of functions that map  $\mathcal{D}_\lambda := \{0, 1\}^\ell$  to a group  $\mathcal{G}_\lambda$ , that supports aggregation over sets recognized by decision trees. That is, let  $\mathcal{S}_\lambda = \{S \subseteq \{0, 1\}^\ell : \exists \text{ a decision tree } T_S \in \mathcal{T}_\lambda \text{ that recognizes } S\}$ .

Our construction uses a hypercube-aggregate PRF family  $\mathcal{F}'_\ell$  as a sub-routine. First, we need the following simple lemma.

**Lemma 3 (Decision Trees as Disjoint Unions of Hypercubes).** *Let  $S \subseteq \{0, 1\}^\ell$  be recognized by a decision tree  $T_S$  of size  $p = p(\lambda)$ . Then,  $S$  is a disjoint union of at most  $p$  hypercubes  $H_{y_1}, \dots, H_{y_p}$ , where each  $y_i \in \{0, 1, \star\}^\ell$  and  $H_{y_i} = \text{Match}(y_i)$ . Furthermore, given  $T_S$ , one can in polynomial time compute these hypercubes.*

Given the lemma, **Aggregate** is simple: on input a set  $S$  represented by a decision tree  $T_S$ , compute the disjoint hypercubes  $H_{y_1}, \dots, H_{y_p}$ . Run the hypercube aggregation algorithm to compute

$$g_i \leftarrow \text{Aggregate}_{\mathcal{F}}(K, y_i)$$

and outputs  $g := \prod_{i=1}^p g_i$ .

Basing the construction on the hypercube-aggregate PRF scheme from Section 3.2, we get a decision tree-aggregate PRF based on the sub-exponential DDH assumption. The security of this PRF follows from Lemma 1 by an argument identical to the one in Section 3.2.

### 3.4 Read-once formulas

Read-once boolean formula provide a different generalization of hypercubes and they too admit an efficient aggregation algorithm for the Naor-Reingold PRF, with a similar security guarantee.

A boolean formula on  $\ell$  variables is a circuit on  $x = (x_1, \dots, x_\ell) \in \{0, 1\}^\ell$  composed of only AND, OR, and NOT gates. A *read-once boolean formula* is a boolean formula with fan-out 1, namely each input literal feeds into at most one gate, and each gate output feeds into at most one other gate.<sup>13</sup> Let  $R_\lambda$  be the family of all read-once boolean formulas over  $\ell(\lambda)$  variables. Without loss of generality, we restrict these circuits to be in a standard form: namely, composed of fan-in 2 and fan-out 1 AND and OR gates, and any NOT gates occurring at the inputs.

In this form, the circuit for any read-once boolean formula can be identified with a labelled binary tree; we identify a formula by the label of its root  $C_\phi$ .

<sup>13</sup> We allow a formula to ignore some inputs variables; this enables the model to express hypercubes directly.

Nodes with zero children are variables or their negation, labelled by  $x_i$  or  $\bar{x}_i$ , while all other nodes have 2 children and represent gates with fan-in 2. For such a node with label  $C$ , its children have labels  $C_L$  and  $C_R$ . Note that each child is itself a read-once boolean formula on fewer inputs, and their inputs are disjoint. Let the gate type of a node  $C$  be  $\text{type}(C) \in \{AND, OR\}$ .

We describe a recursive aggregation algorithm for computing products of PRF values over all accepting inputs for a given read-once boolean formula  $C_\phi$ . Looking forward, we require the formula to be read-once in order for the recursion to be correct. The algorithm described reduces to that of Section 3.2 in the case where  $\phi$  describes a hypercube.

**Construction** The aggregation algorithm for read-once Boolean formulas takes as input the PRF key  $\mathbf{K} = (K_0, \dots, K_\ell)$  and a formula  $C_\phi \in R_\lambda$  where  $C_\phi$  only reads the variables  $x_1, \dots, x_m$  for some  $m \leq \ell$ . We abuse notation and interpret  $C_\phi$  to be a formula on both  $\{0, 1\}^\ell$  and  $\{0, 1\}^m$  in the natural way.

$$\text{AGG}_{k, \Pi}(C_\phi) = \prod_{x: C_\phi(x)=1} g^{K_0 \prod_{i \in [\ell]} K_i^{x_i}} \quad (1)$$

$$= g^{K_0 \sum_{x: C_\phi(x)=1} \prod_{i \in [\ell]} K_i^{x_i}} \quad (2)$$

$$= g^{K_0 \cdot A(C_\phi, 1) \cdot \prod_{m < j \leq \ell} (1 + K_j)} \quad (3)$$

where we define  $A(C, 1) := \sum_{\{x \in \{0, 1\}^m : C(x)=1\}} \prod_{i \in [m]} K_i^{x_i}$ . If  $A(C, 1)$  is efficiently computable, then **Aggregate** will simply compute it and return (3). To this end, we provide a recursive procedure for computing  $A(C, 1)$ .

Generalizing the definition for any sub-formula  $C$  with variables named  $x_1$  to  $x_m$ , define the values  $A(C, 0)$  and  $A(C, 1)$ :

$$A(C, b) := \sum_{\{x \in \{0, 1\}^m : C(x)=b\}} \prod_{i \in [m]} K_i^{x_i}.$$

Recursively compute  $A(C, b)$  as follows:

- If  $C$  is a literal for variable  $x_i$ , then by definition:

$$A(C, b) = \begin{cases} K_i & \text{if } C = x_i \\ 1 & \text{if } C = \bar{x}_i \end{cases}$$

- Else, if  $\text{type}(C) = AND$ : Let  $C_L$  and  $C_R$  be the children of  $C$ . By hypothesis, we can recursively compute  $A(C_L, b)$  and  $A(C_R, b)$  for  $b \in \{0, 1\}$ . Compute  $A(C, b)$  as:

$$A(C, 1) = A(C_L, 1) \cdot A(C_R, 1)$$

$$A(C, 0) = A(C_L, 0) \cdot A(C_R, 0) + A(C_L, 1) \cdot A(C_R, 0) + A(C_L, 0) \cdot A(C_R, 1)$$

- Else,  $\text{type}(C) = OR$ : Let  $C_L$  and  $C_R$  be the children of  $C$ . By hypothesis, we can recursively compute  $A(C_L, b)$  and  $A(C_R, b)$  for  $b \in \{0, 1\}$ . Compute  $A(C, b)$  as:

$$\begin{aligned} A(C, 1) &= A(C_L, 1) \cdot A(C_R, 1) + A(C_L, 1) \cdot A(C_R, 0) + A(C_L, 0) \cdot A(C_R, 1) \\ A(C, 0) &= A(C_L, 0) \cdot A(C_R, 0) \end{aligned}$$

**Lemma 4.**  $A(C, b)$  as computed above is equal to  $\sum_{\{x \in \{0,1\}^m : C(x)=b\}} \prod_{i \in [m]} K_i^{x_i}$

*Proof.* For  $C$  a literal, the correctness is immediate. We must check the recursion for each  $\text{type}(C) \in \{AND, OR\}$  and  $b \in \{0, 1\}$ . We only show the case for  $b = 1$  when  $C$  is an OR gate; the other three cases can be shown similarly.

Let  $S_{b_L, b_R} = \{x = (x_L, x_R) : (C_L(x_L), C_R(x_R)) = (b_L, b_R)\}$  be the set of inputs  $(x_L, x_R)$  to  $C$  such that  $C_L(x_L) = b_L$  and  $C_R(x_R) = b_R$ . The set  $\{x : C(x) = 1\}$  can be decomposed into the disjoint union  $S_{0,1} \sqcup S_{1,0} \sqcup S_{1,1}$ . Furthermore,

$$A(C, 1) = \sum_{x \in S_{0,1}} \prod_{i \in [m]} K_i^{x_i} + \sum_{x \in S_{1,0}} \prod_{i \in [m]} K_i^{x_i} + \sum_{x \in S_{1,1}} \prod_{i \in [m]} K_i^{x_i}$$

Because  $C$  is read-once, the sets of inputs on which  $C_L$  and  $C_R$  depend are disjoint; this implies that  $A(C_L, b_L) \cdot A(C_R, b_R) = \sum_{x \in S_{b_L, b_R}} \prod_{i \in [m]} K_i^{x_i}$ , yielding the desired recursion.

**Theorem 12.** Let  $\epsilon > 0$  be a constant, choose the security parameter  $\lambda = \Omega(\ell^{1/\epsilon})$ , and assume  $(2^{\lambda^\epsilon}, 2^{-\lambda^\epsilon})$ -hardness of the DDH assumption. Then, the collection of functions  $\mathcal{F}_\lambda$  defined above is a secure aggregate PRF with respect to the subsets  $R_\lambda$  and the product aggregation function over the group  $\mathcal{G}$ .

*Proof.* Correctness is immediate from Lemma 4, and Equation (3). Security follows from the decisional Diffie-Hellman assumption in much the same way it did in the case of bit-fixing functions.

## 4 Connection to Learning

We now turn a discussion of the connection between augmented PRFs and computational learning. After the preliminaries, we present the learning with aggregate queries model and the corresponding hardness results implied by our constructions of AGG-PRFs. We defer the discussion of learning with restriction queries and access to related concepts to the full version of this paper.

### 4.1 Preliminaries

**Notation:** For a probability distribution  $D$  over a set  $X$ , we denote by  $x \leftarrow D$  to mean that  $x$  is sampled according to  $D$ , and  $x \leftarrow X$  to denote uniform sampling from  $X$ . For an algorithm  $A$  and a function  $\mathcal{O}$ , we denote that  $A$  has oracle access to  $\mathcal{O}$  by  $A^{\mathcal{O}(\cdot)}$ .

We recall the definition of a “concept class”. In this section, we will often need to explicitly reason about the representations of the concept classes discussed. Therefore we make use of the notion of a “representation class” as defined by [KV94] alongside that of concept classes. This unified formalization enables us to discuss both these traditional learning models (namely, PAC and learning with membership queries) as well as the new models we present below. Our definitions are parametrized by  $\lambda \in \mathbb{N}$ .<sup>14</sup>

**Definition 5 (Representation class [KV94]).** *Let  $K = \{K_\lambda\}_{\lambda \in \mathbb{N}}$  be a family of sets, where each  $k \in K_\lambda$  has description in  $\{0, 1\}^{s_k(\lambda)}$  for some polynomial  $s_k(\cdot)$ . Let  $X = \{X_\lambda\}_{\lambda \in \mathbb{N}}$  be a set, where each  $X_\lambda$  is called a domain and each  $x \in X_\lambda$  has description in  $\{0, 1\}^{s_x(\lambda)}$  for some polynomial  $s_x(\cdot)$ . With each  $\lambda$  and each  $k \in K_\lambda$ , we associate a Boolean function  $f_k : X_\lambda \rightarrow \{0, 1\}$ .<sup>15</sup> We call each such function  $f_k$  a concept, and  $k$  its index or its description. For each  $\lambda$ , we define the concept class  $C_\lambda = \{f_k : k \in K_\lambda\}$  to be the set of all concepts with index in  $K_\lambda$ . We define the representation class  $C = \{C_\lambda\}$  to be the union of all concept classes  $C_\lambda$ .*

This formalization allows us to easily associate complexity classes with concepts in learning theory. For example, to capture the set of all DNF formulas on  $\lambda$  inputs with size at most  $p(\lambda)$  for a polynomial  $p$ , we will let  $X_\lambda = \{0, 1\}^\lambda$ , and  $K_\lambda^{p(\lambda)}$  be the set of descriptions of all DNF formulas on  $\lambda$  variables with size at most  $p(\lambda)$  under some reasonable representation. Then a concept  $f_k(x)$  evaluates the formula  $k$  on input  $x$ . Finally,  $\text{DNF}_\lambda^{p(\lambda)} = \{f_k : k \in K_\lambda^{p(\lambda)}\}$  is the concept class, and  $\text{DNF}^{p(\lambda)} = \{\text{DNF}_\lambda^{p(\lambda)}\}_{\lambda \in \mathbb{N}}$ .  $\text{DNF}^{p(\lambda)}$  is the representation class that computes all DNF formulas on  $\lambda$  variables with description of size at most  $p(\lambda)$  in the given representation.

As a final observation, note that a Boolean-valued PRF family  $\mathcal{F} = \{\mathcal{F}_\lambda\}$  where  $\mathcal{F}_\lambda = \{f_k : X_\lambda \rightarrow \{0, 1\}\}$  with keyspace  $K = \{K_\lambda\}$  and domain  $X = \{X_\lambda\}$  satisfies the syntax of a representation class as defined above. This formalization is useful precisely because it captures both PRF families and complexity classes, enabling lower bounds in various learning models.

In proving lower bounds for learning representation classes, it will be convenient to have a notion of containment for two representation classes.

**Definition 6 ( $\subseteq$ ).** *For two representation classes  $\mathcal{F} = \{\mathcal{F}_\lambda\}$  and  $\mathcal{G} = \{\mathcal{G}_\lambda\}$  on the same domain  $X = \{X_\lambda\}$ , and with indexing sets  $I = \{I_\lambda\}$  and  $K = \{K_\lambda\}$  respectively, we say  $\mathcal{F} \subseteq \mathcal{G}$  if for all sufficiently large  $\lambda$ , for all  $i \in I_\lambda$ , there exists  $k \in K_\lambda$  such that  $g_k \equiv f_i$ .*

Informally, if a representation class contains a PRF family, then this class is hard to MQ-learn (as in [Val84]). We apply similar reasoning to more powerful learning models. For example, if  $\mathcal{G}$  is the representation class  $\text{DNF}^{p(\lambda)}$  as defined

<sup>14</sup> When clear from the context, we will omit the subscript  $\lambda$ .

<sup>15</sup> This association is an efficient procedure for evaluating  $f_k$ . Concretely, we might consider that there is a universal circuit  $F_\lambda$  such that for each  $\lambda$ ,  $f_k(\cdot) = F_\lambda(k, \cdot)$ .

above, then  $\mathcal{F} \subseteq \text{DNF}^{p(\lambda)}$  is equivalent to saying that for all sufficiently large  $\lambda$ , the concept class  $\mathcal{F}_\lambda$  can be decided by a DNF on  $\lambda$  inputs of  $p(\lambda)$  size.

We now recall some standard definitions.

**Definition 7 ( $\epsilon$ -approximation).** *Let  $f, h : X \rightarrow \{0, 1\}$  be arbitrary functions. We say  $h$   $\epsilon$ -approximates  $f$  if  $\Pr_{x \leftarrow X}[h(x) \neq f(x)] \leq \epsilon$ .*

In general,  $\epsilon$ -approximation is considered under a general distribution on  $X$ , but we will consider only the uniform distribution in this work.

**Definition 8 (PAC learning).** *For a concept  $f : X_\lambda \rightarrow \{0, 1\}$ , and a probability distribution  $D_\lambda$  over  $X_\lambda$ , the example oracle  $EX(f, D_\lambda)$  takes no input and returns  $(x, f(x))$  for  $x \leftarrow D_\lambda$ . An algorithm  $\mathcal{A}$  is an  $(\epsilon, \delta)$ -PAC learning algorithm for representation class  $\mathcal{C}$  if for all sufficiently large  $\lambda$ ,  $\epsilon = \epsilon(\lambda) > 0$ ,  $\delta = \delta(\lambda) > 0$  and  $f \in \mathcal{C}_\lambda$ ,*

$$\Pr[\mathcal{A}^{EX(f, D_\lambda)} = h : h \text{ is an } \epsilon\text{-approximation to } f] \geq 1 - \delta$$

**Definition 9 (MQ learning).** *For a concept  $f : X_\lambda \rightarrow \{0, 1\}$ , the membership oracle  $MEM(f)$  takes as input a point  $x \in X_\lambda$  and returns  $f(x)$ . An algorithm  $\mathcal{A}$  is an  $(\epsilon, \delta)$ -MQ learning algorithm for representation class  $\mathcal{C}$  if for all sufficiently large  $\lambda$ ,  $\epsilon = \epsilon(\lambda) > 0$ ,  $\delta = \delta(\lambda) > 0$ , and  $f \in \mathcal{C}_\lambda$ ,*

$$\Pr[\mathcal{A}^{MEM(f)} = h : h \text{ is an } \epsilon\text{-approximation to } f] \geq 1 - \delta$$

We consider only PAC learning *with uniform examples*, where  $D_\lambda$  is the uniform distribution over  $X_\lambda$ . In this case, MQ is strictly stronger than PAC: everything that is PAC learnable is MQ learnable.

Observe that for any  $f : X_\lambda \rightarrow \{0, 1\}$ , either  $h(x) = 0$  or  $h(x) = 1$  will  $\frac{1}{2}$ -approximate  $f$ . Furthermore, if  $\mathcal{A}$  is inefficient,  $f$  may be learned exactly. For a learning algorithm to be non-trivial, we require that it is *efficient* in  $\lambda$ , and that it at least *weakly* learns  $\mathcal{C}$ .

**Definition 10 (Efficient- and weak- learning).**

- $\mathcal{A}$  is said to be *efficient* if the time complexity of  $\mathcal{A}$  and  $h$  are polynomial in  $1/\epsilon, 1/\delta$ , and  $\lambda$ .
- $\mathcal{A}$  is said to *weakly learn*  $\mathcal{C}$  if there exist some polynomials  $p_\epsilon(\lambda), p_\delta(\lambda)$  for which  $\epsilon \leq \frac{1}{2} - \frac{1}{p_\epsilon(\lambda)}$  and  $\delta \leq 1 - \frac{1}{p_\delta(\lambda)}$ .
- We say a representation class is *learnable* if it is both *efficiently* and *weakly learnable*. Otherwise, it is *hard to learn*.

Lastly, we recall the efficiently recognizable ensembles of sets as defined in Section 2. We occasionally call such ensembles indexed, or succinct. Throughout this section, we require this property of our set ensembles  $\mathcal{S}$ . Both the  $\text{MQ}_{\text{RA}}$  and  $\text{AQ}$  learning models that we present are defined with respect to  $\mathcal{S} = \{\mathcal{S}_\lambda\}$ , an efficiently recognizable ensemble of subsets of the domain  $X_\lambda$ .



## 4.2 Learning with Aggregate Queries

This computational learning model is inspired by our aggregate PRFs. Rather than being a natural model in its own right, this model further illustrates how cryptography and learning are in some senses duals. Here, we consider a new extension to the power of the learning algorithm. Whereas membership queries are of the form “What is the label of an example  $x$ ?”, we grant the learner the power to request the evaluation of simple functions on tuples of examples  $(x_1, \dots, x_k)$  such as “How many of  $(x_1 \dots x_k)$  are in  $C$ ?” or “Compute the product of the labels of  $(x_1, \dots, x_k)$ ?”. Clearly, if  $k$  is polynomial then this will result only a polynomial gain in the query complexity of a learning algorithm in the best case. Instead, we propose to study cases when  $k$  may be super polynomial, but the description of the tuples is succinct. For example, the learning algorithm might query the number of  $x$ ’s in a large interval that are positive examples in the concept.

As with the restriction access and related concept models – and the aggregate PRFs we define in this work – the Aggregate Queries (AQ) learning model will be considered with restrictions to both the types of aggregate functions  $\Gamma$  the learner can query, and the sets  $\mathcal{S}$  over which the learner may request these functions to be evaluated on. We now present the AQ learning model informally:

**Definition 11 (( $\Gamma, \mathcal{S}$ )-Aggregate Queries (AQ) Learning).** *Let  $\mathcal{C}$  be a representation class with domains  $X = \{X_\lambda\}$ , and  $\mathcal{S} = \{\mathcal{S}_\lambda\}$  where each  $\mathcal{S}_\lambda$  is a collection of efficiently recognizable subsets of the  $X_\lambda$ .  $\Gamma : \{0, 1\}^* \rightarrow V_\lambda$  be an aggregation function [as in def:]. Let  $\text{AGG}_k^\lambda \triangleq \text{AGG}_{f_k, \mathcal{S}_\lambda, \Gamma_\lambda}^\lambda$  be the aggregation oracle for  $f_k \in \mathcal{C}_\lambda$ , for  $S \in \mathcal{S}_\lambda$  and  $\Gamma_\lambda$ .*

*An algorithm  $\mathcal{A}$  is an  $(\epsilon, \delta)$ -( $\Gamma, \mathcal{S}$ )-AQ learning algorithm for  $\mathcal{C}$  if, for all sufficiently large  $\lambda$ , for every  $f_k \in \mathcal{C}_\lambda$ ,  $\Pr[\mathcal{A}^{\text{MEM}_{f_k}(\cdot), \text{AGG}_{f_k}^\lambda(\cdot)} = h] \geq 1 - \delta$  where  $h$  is an  $\epsilon$ -approximation to  $f_k$ .*

### Hardness of aggregate query learning

**Theorem 13.** *Let  $\mathcal{F}$  be a boolean-valued aggregate PRF with respect to set system  $\mathcal{S} = \{\mathcal{S}_\lambda\}$  and accumulation function  $\Gamma = \{\Gamma_\lambda\}$ . For a representation class  $\mathcal{C}$ , if  $\mathcal{F} \subseteq \mathcal{C}$ , then  $\mathcal{C}$  is hard to  $(\Gamma, \mathcal{S})$ -AQ learn.*

Looking back to our constructions of aggregate pseudorandom function families from the prequel, we have the following corollaries.

**Corollary 4.** *The existence of one-way functions implies that  $P/\text{poly}$  is hard to  $(\sum, \mathcal{S}_{[a,b]})$ -AQ learn, with  $\mathcal{S}_{[a,b]}$  the set of sub-intervals of the domain as defined in Section 3.*

**Corollary 5.** *The DDH Assumption implies that  $NC^1$  is hard to  $(\sum, \mathcal{S}_{[a,b]})$ -AQ learn, with  $\mathcal{S}_{[a,b]}$  the set of sub-intervals of the domain as defined in Section 3.*

**Corollary 6.** *The subexponential DDH Assumption implies that  $NC^1$  is hard to  $(\prod, \mathcal{R})$ -AQ learn, with  $\mathcal{R}$  the set of read-once boolean formulas defined in Section 3.*

*Proof (Proof of Theorem 13).* Interpreting  $\mathcal{F}$  itself as a concept class, we will show an efficient reduction from violating the pseudorandomness property of  $\mathcal{F}$  to weakly  $(\Gamma, \mathcal{S})$ -AQ learning  $\mathcal{F}$ . By assumption,  $\mathcal{F} \subseteq \mathcal{C}$ , implying that  $\mathcal{C}$  is hard to learn as well.

**Reduction:** Suppose for contradiction that there exists an efficient weak learning algorithm  $\mathcal{A}$  for  $\mathcal{F}$ . We define algorithm  $\mathcal{B}$  violating the aggregate PRF security of  $\mathcal{F}$ . In the PRF security game,  $\mathcal{B}$  is presented with two oracles:  $F(\cdot)$  and  $AGG_F^\lambda$  for a function  $F$  chosen according to the secret bit  $b \in \{0, 1\}$ . In EXP(0),  $F = f_k$  for random  $k \in K_\lambda$ ; by assumption  $f_k \in \mathcal{C}_\lambda$ . In EXP(1),  $F$  is a uniformly random function from  $X$  to  $\{0, 1\}$ . The learning algorithm  $\mathcal{A}$  is presented with precisely the same oracles.  $\mathcal{B}$  runs  $\mathcal{A}$ , simulating its oracles by passing queries and responses to its own oracles.  $X_{\mathcal{A}} = \{x \in X_\lambda : \mathcal{A} \text{ queried } (\psi, x) \text{ for some } \psi\}$ . Once  $\mathcal{A}$  terminates, it outputs hypothesis  $h$ .

After receiving hypothesis  $h$ ,  $\mathcal{B}$  estimates the probability

$$p = \Pr_{x \leftarrow X \setminus X_{\mathcal{A}}} [h(x) = F(x)]$$

(using polynomial in  $\lambda, p_\epsilon(\lambda)$  samples). In EXP(0), this probability is at least  $1 - \epsilon$  with probability at least  $1 - \delta$ ; in EXP(1), it is exactly  $1/2$ . To sample uniform  $x \in X \setminus X_{\mathcal{A}}$ , we simply take a uniform  $x \in X$ : with high probability  $x \in X \setminus X_{\mathcal{A}}$ . If the estimate is close to  $\epsilon$ , guess EXP(0); otherwise, flip an fair coin  $b' \in \{0, 1\}$  and guess EXP( $b'$ ). The advantage  $ADV_\lambda^{APRF}$  of  $\mathcal{B}$  in the PRF security game is at least  $\frac{1}{3p_\delta(n)}$  for all sufficiently large  $\lambda$  (as shown below), directly violating the security of  $\mathcal{F}$ .

Let

$$p_b \triangleq \Pr_{x \in X \setminus X_{\mathcal{A}}} [h(x) \neq F(x) | EXP(b)]$$

be the probability taken with respect to experiment EXP( $b$ ). In EXP(1),  $F$  is a uniformly random function. Thus,  $p_1 = \frac{1}{2}$ . With high probability,  $\mathcal{B}$  will output a random bit  $b' \in \{0, 1\}$ , guessing correctly with probability  $1/2$ .

In EXP(0),  $h$  is an  $\epsilon$ -approximation to  $F$  with probability at least  $1 - \delta$ . In this case,  $p_0 \geq 1 - \epsilon \geq \frac{1}{2} + \frac{1}{p_\epsilon(\lambda)}$ . By a Hoeffding bound,  $\mathcal{B}$  will guess  $b' = 0$  with high probability by estimating  $p$  using only polynomial in  $\lambda, p_\epsilon(\lambda)$  samples. On the other hand, if  $h$  is not an  $\epsilon$ -approximation,  $\mathcal{B}$  will  $b' = 0$  with probability at least  $1/2$ .

Let  $\text{negl}(\lambda)$  be the error probability from the Hoeffding bound, which can be made exponentially small in  $\lambda$ . The success probability is:

$$\Pr[b = b' | b = 0] \geq (1 - \delta)(1 - \text{negl}(\lambda)) + \frac{\delta}{2}$$

which, for  $1 - \delta \geq \frac{1}{p_\delta(\lambda)}$  is at least  $\frac{1}{3p_\delta(\lambda)} + \frac{1}{2}$  for sufficiently large  $\lambda$ . Thus  $\mathcal{B}$  a non-negligible advantage of  $1/3p_\delta(\lambda)$  in the  $(\Gamma, \mathcal{S})$ -aggregate-PRF security game.

**Acknowledgements** Aloni Cohen’s research was supported in part by NSF Graduate Student Fellowship and NSF grants CNS1347364, CNS1413920 and FA875011-20225.

Shafi Goldwasser’s research was supported in part by NSF Grants CNS1347364, CNS1413920 and FA875011-20225.

Vinod Vaikuntanathan’s research was supported by DARPA Grant number FA8750-11-2-0225, an Alfred P. Sloan Research Fellowship, an NSF CAREER Award CNS-1350619, NSF Frontier Grant CNS-1414119, a Microsoft Faculty Fellowship, and a Steven and Renee Finn Career Development Chair from MIT.

## References

- ABPP14. Michel Abdalla, Fabrice Benhamouda, Alain Passelègue, and Kenneth G. Paterson. Related-key security for pseudorandom functions beyond the linear barrier. 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 77–94. Springer, 2014.
- BC10. Mihir Bellare and David Cash. Pseudorandom functions and permutations provably secure against related-key attacks. In Tal Rabin, editor, *Advances in Cryptology - CRYPTO 2010, 30th Annual Cryptology Conference, Santa Barbara, CA, USA, August 15-19, 2010. Proceedings*, volume 6223 of *Lecture Notes in Computer Science*, pages 666–684. Springer, 2010.
- BF02. Nader H Bshouty and Vitaly Feldman. On using extended statistical queries to avoid membership queries. *The Journal of Machine Learning Research*, 2:359–395, 2002.
- BGI14a. Elette Boyle, Shafi Goldwasser, and Ioana Ivan. Functional signatures and pseudorandom functions. In Krawczyk [Kra14], pages 501–519.
- BGI14b. Elette Boyle, Shafi Goldwasser, and Ioana Ivan. Functional signatures and pseudorandom functions. In Krawczyk [Kra14], pages 501–519.
- BGV11. Siavosh Benabbas, Rosario Gennaro, and Yevgeniy Vahlis. Verifiable delegation of computation over large datasets. In Phillip Rogaway, editor, *Advances in Cryptology - CRYPTO 2011 - 31st Annual Cryptology Conference, Santa Barbara, CA, USA, August 14-18, 2011. Proceedings*, volume 6841 of *Lecture Notes in Computer Science*, pages 111–131. Springer, 2011.
- BLMR13. Dan Boneh, Kevin Lewi, Hart William Montgomery, and Ananth Raghunathan. Key homomorphic prfs and their applications. In Ran Canetti and Juan A. Garay, editors, *Advances in Cryptology - CRYPTO 2013 - 33rd Annual Cryptology Conference, Santa Barbara, CA, USA, August 18-22, 2013. Proceedings, Part I*, volume 8042 of *Lecture Notes in Computer Science*, pages 410–428. Springer, 2013.
- BW13a. Dan Boneh and Brent Waters. Constrained pseudorandom functions and their applications. In Sako and Sarkar [SS13], pages 280–300.
- BW13b. Dan Boneh and Brent Waters. Constrained pseudorandom functions and their applications. In Sako and Sarkar [SS13], pages 280–300.
- DRWY12. Zeev Dvir, Anup Rao, Avi Wigderson, and Amir Yehudayoff. Restriction access. In Shafi Goldwasser, editor, *Innovations in Theoretical Computer*

- Science 2012, Cambridge, MA, USA, January 8-10, 2012*, pages 19–33. ACM, 2012.
- GGM86. Oded Goldreich, Shafi Goldwasser, and Silvio Micali. How to construct random functions. *J. ACM*, 33(4):792–807, 1986. Extended abstract in FOCS 84.
- GGN10. Oded Goldreich, Shafi Goldwasser, and Asaf Nussboim. On the implementation of huge random objects. *SIAM J. Comput.*, 39(7):2761–2822, 2010.
- HILL99. Johan Håstad, Russell Impagliazzo, Leonid A. Levin, and Michael Luby. A pseudorandom generator from any one-way function. *SIAM J. Comput.*, 28(4):1364–1396, 1999.
- KPTZ13a. Aggelos Kiayias, Stavros Papadopoulos, Nikos Triandopoulos, and Thomas Zacharias. Delegatable pseudorandom functions and applications. In Sadeghi et al. [SGY13], pages 669–684.
- KPTZ13b. Aggelos Kiayias, Stavros Papadopoulos, Nikos Triandopoulos, and Thomas Zacharias. Delegatable pseudorandom functions and applications. In Sadeghi et al. [SGY13], pages 669–684.
- Kra14. Hugo Krawczyk, editor. *Public-Key Cryptography - PKC 2014 - 17th International Conference on Practice and Theory in Public-Key Cryptography, Buenos Aires, Argentina, March 26-28, 2014. Proceedings*, volume 8383 of *Lecture Notes in Computer Science*. Springer, 2014.
- KV94. Michael J. Kearns and Leslie G. Valiant. Cryptographic limitations on learning boolean formulae and finite automata. *J. ACM*, 41(1):67–95, 1994.
- NR04. Moni Naor and Omer Reingold. Number-theoretic constructions of efficient pseudo-random functions. *J. ACM*, 51(2):231–262, 2004.
- RR97. Alexander A. Razborov and Steven Rudich. Natural proofs. *J. Comput. Syst. Sci.*, 55(1):24–35, 1997.
- SGY13. 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*. ACM, 2013.
- SS13. Kazue Sako and Palash Sarkar, editors. *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*, volume 8270 of *Lecture Notes in Computer Science*. Springer, 2013.
- 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.
- Val84. Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- VV86. Leslie G. Valiant and Vijay V. Vazirani. NP is as easy as detecting unique solutions. *Theor. Comput. Sci.*, 47(3):85–93, 1986.