

Combiners for Backdoored Random Oracles

Balthazar Bauer^{2,1}, Pooya Farshim^{1,2}, and Sogol Mazaheri³

¹ DI/ENS, CNRS, PSL University, Paris, France

² Inria, Paris, France

balthazar.bauer@ens.fr, pooya.farshim@gmail.com

³ Cryptoplexity, Technische Universität Darmstadt, Germany

sogol.mazaheri@cryptoplexity.de

Abstract. We formulate and study the security of cryptographic hash functions in the *backdoored random-oracle* (BRO) model, whereby a big brother designs a “good” hash function, but can also see *arbitrary functions* of its table via backdoor capabilities. This model captures intentional (and unintentional) weaknesses due to the existence of collision-finding or inversion algorithms, but goes well beyond them by allowing, for example, to search for structured preimages. The latter can easily break constructions that are secure under random inversions.

BROs make the task of bootstrapping cryptographic hardness somewhat challenging. Indeed, with only a single arbitrarily backdoored function no hardness can be bootstrapped as any construction can be inverted. However, when two (or more) independent hash functions are available, hardness emerges *even with unrestricted and adaptive access to all backdoor oracles*. At the core of our results lie new reductions from cryptographic problems to the *communication complexities* of various two-party tasks. Along the way we establish a communication complexity lower bound for set-intersection for cryptographically relevant ranges of parameters and distributions and where set-disjointness can be easy.

Keywords. Random oracle, combiner, communication complexity, set-disjointness, set-intersection, lower bounds.

1 Introduction

Hash functions are one of the most fundamental building blocks in the design of cryptographic protocols. From a provable security perspective, a particularly successful methodology to use hash functions in protocols has been the introduction of the random-oracle (RO) model [15,5]. This model formalizes the intuition that the outputs of a well-designed hash function look random by giving all parties, honest or otherwise, oracle access to a uniformly chosen random function. The strong randomness properties inherent in the oracle, in turn, facilitate the security analyses of many protocols.

The cryptanalytic validation of hash functions can strengthen our confidence in this RO-like behavior. On the other hand, as such analyses improve, (unintentional) weaknesses in hash functions are discovered, which can lead to their partial or total break of security. However, cryptanalytic validation might also fail to detect *intentional* weaknesses that are built into systems. For example

such backdoors might be themselves built using cryptographic techniques, which make them hard to detect. Prominent examples show that such backdoors exist and can be exploited in various ways [6,10,11].

In this work we revisit a classical question on protecting against failures of hash functions. Numerous works in this area have studied if, and to what level, by *combining* different hash functions one can offer such protections; see [7,16,17,20] for theoretical treatments and [30,26,13] for cryptanalytic work. However, most work has their focus on unintentional failures (to protect against cryptanalytic advances). In this work, we consider a more adversarial view of hash function failures and ask if well-designed, but possibly *backdoored* hash functions can be used to build backdoor-free hash functions?

Depending on what well-designed means, what adversarial powers the backdoors provide, and what security goals are targeted, different solutions emerge. Hash-function combiners in the works above typically convert two or more hash functions into a new one that is secure as long as *any* of the underlying hash functions is secure. For example, the concatenation combiner builds a collision-resistant hash function given k hash functions as long as one function is collision resistant. Multi-property combiners for other notions, such as PRG, MAC or PRF security, also exist [17].

Typical combiners, however, do not necessarily offer protection when *all* hash functions fail. Intuitively, the goal here is more challenging as all “sources of hardness” have been rendered useless. Despite this, a number of works [27,20,26,33,23] take a more practical approach and introduce an intermediate *weakened* RO model, where hash functions are vulnerable to strong forms of attack, but are otherwise random.

This is an approach that we also adopt here. Since our goal is to protect against *adversarial* weaknesses (aka. backdoors), we place no assumptions on hash-function weaknesses—they can go well beyond computing random preimages or collisions.

1.1 Contributions

We introduce a substantially weakened RO model where an adversary, on top of hash values, can also obtain *arbitrary functions* of the table of the hash function. We formalize this capability via access to a *backdoor oracle* $\text{BD}(f)$ that on input a function f returns $f(\langle \mathbf{H} \rangle)$, arbitrary auxiliary information about the function table of the hash function \mathbf{H} . We call this the *backdoored random-oracle* (BRO) model.

Such backdoors are powerful enough to allow for point inversions—simply hardwire the point y that needs to be inverted into a function $f[y]$ that searches for a preimage of y under \mathbf{H} —or finding collisions. But they can go well beyond them. For example, although Liskov [27] proves one-way security of the combiner $\mathbf{H}(0|x_1|x_2)|\mathbf{H}(1|x_2|x_1)$ under *random* inversions, it becomes insecure when inverted points are not assumed to be random: given $y_1|y_2$ simply look for an inverse $0|x'_1|x'_2$ for y_1 such that $1|x'_2|x'_1$ also maps to y_2 . BRO can also model

arbitrary *preprocessing* attacks (aka. non-uniform attacks) as any auxiliary information about $\langle H \rangle$ can be computed via a one-time oracle access at the onset. This means that collisions (without salting) can be easily found. Furthermore, since BD calls can be adaptive, salting does not help in our setting at all. Indeed, with a single hash function and arbitrary backdoor capabilities no combiner can exist as any construction $C^H(x)$ can be easily inverted by a function that sees the entire $\langle H \rangle$ and searches for inversions.

In practice it is natural to assume that independent hash functions are available. We can easily model this by an extension to the k -BRO model, whereby k independent ROs and their respective backdoor oracles are made available.¹ The interpretation in our setting is that different “trusted” authorities have designed and made public hash functions that display good (i.e., RO-like) behaviors, but their respective backdoors enable computing any function of the hash tables. We ask if these hash functions can be combined in way that renders their backdoors useless. We observe that the result of Hoch and Shamir [20] can be seen as one building a collision-resistant hash function in the 2-BRO model assuming backdoor oracles that allow for random inversions only.

From a high-level point of view, our main result shows that in the 2-BRO model cryptographic hardness *can* be bootstrapped, even with access to *both* backdoor oracles and even when *arbitrary* backdoor capabilities are provided. In other words, there are secure constructions in the 2-BRO model that can tolerate arbitrary weaknesses in all underlying hash functions. At the core of our results lies new links with hard problems in the area of *communication complexity*.

COMMUNICATION COMPLEXITY. The communication complexity [38,24] of a two-party task $f(S, T)$ is the minimum communication cost over two-party protocols that compute $f(S, T)$. Two rich and well-studied problems in this area are the *set-disjointness* and *set-intersection* problems (see [9] for a survey). Here two parties hold sets S and T respectively. In set-disjointness, their goal is to decide whether or not $S \cap T = \emptyset$; in set-intersection they need to compute at least one element in this intersection. Typically, work in communication complexity studies communication cost over all inputs, that is, the *worst-case* communication complexity of a problem, as the focus is on lower bounds. Cryptographic applications, on the other hand, usually require *average-case* hardness. Distributional (average-case) communication complexity of a problem averages the communication cost over random choices of (S, T) from some distribution μ . We will rely on average-case lower bounds in this work.

THE BASIC IDEAS. In this work, we focus on the parallel (concatenation) and sequential (cascade) composition of hash functions H_1 and H_2 and consider the combiners:

$$C_{\perp}^{H_1, H_2}(x) := H_1(x) \parallel H_2(x) \quad \text{and} \quad C_{\circ}^{H_1, H_2}(x) := H_2(H_1(x)) .$$

¹ Note that k -BRO can be viewed as a *restricted* version of the 1-BRO model where k ROs are built from a single RO acting on k separate domains and backdoor capabilities are restricted to these domains only.

Here $H_1 : \{0, 1\}^n \rightarrow \{0, 1\}^{n+s_1}$ and $H_2 : \{0, 1\}^n \rightarrow \{0, 1\}^{n+s_2}$ in the first construction, and $H_1 : \{0, 1\}^n \rightarrow \{0, 1\}^{n+s_1}$ and $H_2 : \{0, 1\}^{n+s_1} \rightarrow \{0, 1\}^{n+s_1+s_2}$ in the second.

Consider the one-way security of the concatenation combiner in the 2-BRO model. An adversary is given a point $y^* := y_1^* | y_2^* := H_1(x^*) | H_2(x^*)$ for a random x^* . It has access to the backdoor oracles BD_1 and BD_2 for functions H_1 and H_2 respectively. Its goal is to compute a preimage x for y^* under $C_{|}^{H_1, H_2}$. This is the case iff $H_1(x) = y_1^*$ and $H_2(x) = y_2^*$. Now define two sets $S := H_1^{-1}(y_1^*)$, the set of preimages of y_1^* under H_1 , and $T := H_2^{-1}(y_2^*)$, the set of preimages of y_2^* under H_2 . *Thus the adversary wins iff $x \in S \cap T$.*

The two backdoor oracles respectively know S and T as they are part of the descriptions of the two hash functions. This allows us to convert a successful one-way adversary to a two-party protocol that computes an element x of the intersection $S \cap T$. Put differently, if the communication complexity of set-intersection for sets that are distributed as above has a high lower bound, then the adversary has to place a large number of queries, which, in turn, allows us to conclude that the concatenation combiner is one-way in the 2-BRO model.

The question is: for which sets S and T is set-intersection hard? Suppose the hash functions $H_1, H_2 : \{0, 1\}^n \rightarrow \{0, 1\}^m$ are compressing and $m = n - s$. Then on average the sets S and T would each have 2^s elements. We can of course communicate these sets in $\mathcal{O}(2^s)$ bits and find a preimage. However, the cost of this attack when s is linear in n (or even super-logarithmic in n) becomes prohibitive. This raises the question if set-intersection is hard for, say, $s = n/2$ and where the distribution over (S, T) is induced by the two hash functions, where except a single element in common (guaranteed to exist by the rules of the one-way game) all others are sampled uniformly and independently at random and included in the sets.

We observe that hardness of the set-disjointness problem implies hardness of set-intersection as the parties can verify that a given element is indeed in both their sets.² Set-disjointness is a better studied problem. To the best of our knowledge two results on set-disjointness with parameters and distributions close to those in our setting have been proven. First, a classical (and technical) result of Babai, Simon and Frankl [1] which shows an $\Omega(\sqrt{N})$ lower bound for random and independent sets S and T of size *exactly* \sqrt{N} in a universe of size N . Second, a result based on information-theoretic arguments due to Bar-Yossef et al. [2], for dependent sets S and T , which has been adapted to *Bernoulli product distributions* in lectures by Moshkovitz and Barak [32, Lecture 9] and Guruswami and Cheraghchi [19, Lecture 21]. The distribution is as follows: for each of the N elements in the universe, independent $\text{Ber}(1/\sqrt{N})$ bits are sampled. (The probability of 1 is $1/\sqrt{N}$.) The sets then consist of all elements for which the bit is set to 1.³ The authors again prove an $\Omega(\sqrt{N})$ lower bound (which

² On the other hand, for sufficiently large sets that intersect with high probability, set-disjointness is easy whereas set-intersection can remain hard.

³ The expected size of such Bernoulli sets is $N/\sqrt{N} = \sqrt{N}$, but this size can deviate from the mean and this distribution is *not* identical to that by Babai et al. [1].

is tight up to logarithmic factors). We note that both these results only hold for protocols that err with probability at most $\varepsilon \leq 1/100$. However, we only found incomplete proofs of set-disjointness for product Bernoulli distributions, and thus have included a self-contained proof in the full version of this paper [4, Appendix C]. We also prove a distributional communication complexity lower bound for set-*intersection* for parameters where set-disjointness can be *easy*.

The second result is better suited for our purposes as the size restriction in the first one would restrict us to regular random oracles. Indeed, the distribution induced on the preimages of y_1^* (resp. y_2^*) by the hash function outside the common random point *is* Bernoulli: $\Pr[\mathbf{H}_1(x) = y_1^*] = 1/2^m$ (resp. $\Pr[\mathbf{H}_2(x) = y_2^*] = 1/2^m$) for any x and independently for values of x . We use this fact to show that set-*intersection* and set-disjointness problems are, respectively, sufficient to prove it is hard to invert random co-domain points (a property that we call random preimage resistance, rPre) or even decide if a preimage exists (which we call oblivious PRG, oPRG). The main benefit of these games is that they do away with the common point guaranteed to exist by the rules of one-way game (and also similar technicalities associated with the standard PRG game). These games can then be related to the one-way and PRG games via cryptographic reductions.

Our lower bound for set-*intersection* allows us to prove strong one-way security for some parameters, while the set-disjointness bound only enables proving weak PRG security. Using amplification techniques we can then convert the weak results to strong one-way functions [18] or strong PRGs [29]. Note that the reductions for all these results are fully black-box and thus would relativize [34]. This implies that the same proofs also hold in the presence of backdoor oracles. Construction of other primitives in minicrypt also relativize. This means we also obtain backdoor-free PRFs, MACs, PRPs, and symmetric encryption schemes in our model. The resulting constructions, however, are often too inefficient to be of any practical use. The bottleneck for PRG efficiency here is the proven lower bounds for set-disjointness. New lower bounds that give trade-offs between protocol error and communication complexity will enable more efficient/secure constructions. We discuss in Section 4 why the current proof does not permit this.

Recall that collision resistance can *not* be based on one-way functions [36]. The concatenation combiner, on the other hand, appears to be collision resistant as *simultaneous collisions* seem hard to find, even with respect to arbitrary backdoors for each hash function. Indeed, an analysis of collision resistance for this combiner reveals a natural *multi-instance* analogue of the set-*intersection* problem, which to the best of our knowledge has not been studied yet. Assuming the hardness of this problem (which we leave open) we get collision resistance. We note that fully black-box amplification for collision-resistance also exists [8], and it is sufficient to prove hardness for small values of protocol error ε (should this be the case as in the case single-instance set-disjointness).

We carry out similar analyses for the *cascade* combiner, for which different choices of parameters lead to security. Although the overall approach remains

Combiner	Strong OW	Weak PRG	Strong CR
Concatenation	$s_1, s_2 = -(\epsilon + 1) \cdot n/2$ for $0 < \epsilon < 1/3$	$s_1 = -n/2 + 1,$ $s_2 = -n/2$	$s_1, s_2 \leq -n/2 - 1$
Cascade	$s_1 = (1 + \epsilon) \cdot n, s_2 = -n$ for $-1/2 < \epsilon < 0$	$s_1 = 2n, s_2 = -2n + 1$	$s_1 = 2n, s_2 = -2n - 1$

Table 1: Overview of results for concatenation and cascade. Functions H_i have stretch s_i . The parameters for collision resistance are conjectural.

the same, we need to deal with difficulties arising from one of the sets being the *image* of a hash function. The latter distribution is somewhat different to Bernoulli sets (as elements are not chosen independently). We show, however, that by addition of noise one-way and PRG security can be based on *known* lower bounds. For collision-resistance we give a reduction to a multi-instance analogue of set-intersection (whose hardness remains open). We analyze the security of the XOR combiner in the full version of this paper [4].

We summarize our results in Table 1. Roughly speaking, strong security demands that the advantage of adversaries in the corresponding security game is negligible, while for weak security it suffices that the advantage is not overwhelming. In the table, concatenation is with respect to hash function $H_1 : \{0, 1\}^n \rightarrow \{0, 1\}^{n+s_1}$ and $H_2 : \{0, 1\}^n \rightarrow \{0, 1\}^{n+s_2}$, while cascade is with respect to hash function $H_1 : \{0, 1\}^n \rightarrow \{0, 1\}^{n+s_1}$ and $H_2 : \{0, 1\}^{n+s_1} \rightarrow \{0, 1\}^{n+s_1+s_2}$. The stretch values s_1 and s_2 can assume negative values (compressing), positive values (expanding), or be zero (length-preserving).

1.2 Discussion

BACKDOORS AS WEAKNESSES. One of the main motivations for the works of Liskov [27] and Hoch and Shamir [20] is the study of design principles for symmetric schemes that can offer protections against weaknesses in their underlying primitives. For example, Hoch and Shamir study the failure-friendly double-pipe hash construction of Lucks [28]. Similarly, Liskov shows that his *zipper hash* is indistinguishable from a random oracle even with an inversion oracle for its underlying compression function. Proofs of security in the unrestricted BRO model would strengthen these results as they place weaker assumptions on the types of weaknesses that are discovered.

AUXILIARY INPUTS. As mentioned above, a closely related model to BRO is the Auxiliary-Input RO (AI-RO) model, introduced by Unruh [37] and recently refined by Dodis, Guo, and Katz [14] and Coretti et al. [12]. Here the result of a one-time preprocessing attack with access to the full table of the random oracle is made available to an adversary. The BRO and AI-RO models are similar in that they both allow for arbitrary functions of the random oracle to be computed. However, BRO allows for adaptive, instance-dependent auxiliary information, whereas the AI-RO model only permits a one-time access at the onset.⁴ Thus

⁴ Arguably, the AI-RO model is better named the Non-Uniform RO model: auxiliary input is often instance dependent whereas non-uniform input is not.

AI-RO is identical to BRO when only a single BD query at the onset is allowed. Extension to multiple ROs can also be considered for AI-ROs, where independent preprocessing attacks are performed on the hash functions. A corollary is that any positive result in the k -BRO model would also hold in the k -AI-RO model. Results in k -AI-RO model can be proven more directly using the decomposition of high-entropy densities as the setting is non-interactive.

FEASIBILITY IN 1-BRO. As already observed, any combiner in 1-BRO is insecure with respect to arbitrary backdoors. We can, however, consider a model where backdoor capabilities are restricted to inversions only. Security in such models will depend on the exact specification of backdoor functionalities \mathcal{F} . For example, under random inversions positive results can be established using standard lazy sampling techniques. But another natural choice is to consider functions which output possibly adversarial preimages, i.e., functions $f[y]$ whose outputs are restricted to those x for which $H(x) = y$. As we have seen, under such generalized inversions provably secure constructions can fail. Moreover, proving security under general inversions seems to require techniques from communication complexity as we do here.

OTHER SETTINGS. Proofs in the random-oracle model often proceed via direct information-theoretic analyses. Here we give cryptographic reductions (somewhat similarly to the standard model) that isolate the underlying communication complexity problems. These problems have diverse applications in other fields (such as circuit complexity, VLSI design, and combinatorial auctions), which motivate their study outside cryptographic contexts. Any improvement in lower bounds for them would also lead to improvements in the security/efficiency of cryptographic constructions. We discussed the benefits of proofs for arbitrary error above. As other examples, results in multi-party communication complexity would translate to the k -BRO model for $k > 2$ or those in quantum communication complexity can be used to built quantum-secure BRO combiners.

1.3 Future work

Our work leaves a number of problems open, some of which are closer to work in communication complexity. We discuss these below.

Lower bounds for set-disjointness that do not assume a small error would improve the security and/or efficiency of our PRG constructions. Moreover, we do not currently have a lower bound for the multi-instance analogue of set-intersection that we need for proving collision resistance. Finding the “maximal” backdoor capabilities in the 1-BRO model under which hardness can be bootstrapped remains an interesting open problem. Katz, Lucks, and Thiruvengadam [22] study the construction of collision-resistant hash functions from ideal ciphers that are vulnerable to differential related-key attacks. We leave the study of combiners for other backdoored primitives, such as ideal ciphers, for future work.

2 Preliminaries

We let \mathbb{N} denote the set of non-negative integers and $\{0, 1\}^n$ be the set of all binary strings of length $n \in \mathbb{N}$. For two bit strings x and y , we denote their concatenation by $x|y$. We let $[N]$ denote the set $\{1, \dots, N\}$. For a finite set S , we denote by $s \leftarrow S$ the uniform random variable over S . The Bernoulli random variable $x \leftarrow \text{Ber}(p)$ takes value 1 with probability p and 0 with probability $1-p$. The Binomial random variable $x_1, \dots, x_n \leftarrow \text{Bin}(n, p)$ constitutes a sequence of n independent Bernoulli samples. We will sometimes use $e^{-x} := \lim_{n \rightarrow \infty} (1-x/n)^n$.

2.1 Random oracles

A hash function H with n -bit inputs and m -bit outputs is simply a function with signature $H : \{0, 1\}^n \rightarrow \{0, 1\}^m$. We let $\text{Fun}[n, m]$ denote the set of all such functions. $\text{Fun}[n, m]$ is finite and we endow it with the uniform distribution. For a hash function H , we let $\langle H \rangle$ denote the function table of H encoded as a string of length $m2^n$. We see the x -th m -bit block of $\langle H \rangle$ as $H(x)$, identifying strings $x \in \{0, 1\}^n$ with integers in $[1, 2^n]$. The random-oracle (RO) model (for a given n and m) is a model of computation where all parties have oracle access to a function $H \leftarrow \text{Fun}[n, m]$.

BACKDOOR FUNCTIONS. A backdoor function for $H \in \text{Fun}[n, m]$ is a function $f : \text{Fun}[n, m] \rightarrow \{0, 1\}^t$. A backdoor capability class \mathcal{F} is a set of such backdoor functions. The unrestricted class contains all functions. But the class can be also restricted, for example, functions $f[y]$ for $y \in \{0, 1\}^m$ whose outputs x are restricted to be in $H^{-}(y)$, where $H^{-}(y)$ is the set preimages of y under H . Randomness can also be hardwired.

THE BRO MODEL. In the backdoored random-oracle (BRO) model, a random function $H \leftarrow \text{Fun}[n, m]$ is sampled. All parties are provided with oracle access to H . Adversarial parties are additionally given access to the procedure

Proc. BD(f) : **return** $f(\langle H \rangle)$

for $f \in \mathcal{F}$. Formally, we denote this model by $\text{BRO}[n, m, \mathcal{F}]$, but will omit $[n, m, \mathcal{F}]$ when it is clear from the context. When $\mathcal{F} = \emptyset$, we recover the conventional RO model. As discussed in the introduction, when the adversarial parties call the backdoor oracle only once and before any hash queries, we recover random oracles with auxiliary input, the AI-RO model [12, Definition 2]. Thus, BRO also models oracle-dependent auxiliary input or pre-computation attacks as special cases. In the k -BRO model (with the implicit parameters $[n_i, m_i, \mathcal{F}_i]$ for $i = 1, \dots, k$) access to k independent random oracles $H_i \in \text{Fun}[n_i, m_i]$ and their respective backdoors BD_i with capabilities \mathcal{F}_i are provided. That is, procedure $\text{BD}_i(f)$ returns $f(\langle H_i \rangle)$. In this work we are primarily interested in the 1-BRO and 2-BRO models with *unrestricted* \mathcal{F} .

We observe that the 2-BRO $[n, m, \mathcal{F}_1, n, m, \mathcal{F}_2]$ model is identical to the 1-BRO $[n+1, m, \mathcal{F}]$ model where for $H \in \text{Fun}[n+1, m]$ we define $H_1(x) := H(0|x)$,

Game $\text{OW}_{\mathcal{C}}^{\mathcal{A}}$	Game $\text{PRG}_{\mathcal{C}}^{\mathcal{A}}$	Game $\text{CR}_{\mathcal{C}}^{\mathcal{A}}$
for $i = 1, 2$ do $\text{H}_i \leftarrow \text{Fun}[n_i, m_i]$ $x \leftarrow \{0, 1\}^n; y \leftarrow \text{C}^{\text{H}_i}(x)$ $x' \leftarrow \mathcal{A}^{\text{H}_i, \text{BD}_i}(y)$ return $(\text{C}^{\text{H}_i}(x') = y)$	for $i = 1, 2$ do $\text{H}_i \leftarrow \text{Fun}[n_i, m_i]$ $y_0 \leftarrow \{0, 1\}^m; b \leftarrow \{0, 1\}$ $x \leftarrow \{0, 1\}^n; y_1 \leftarrow \text{C}^{\text{H}_i}(x)$ $b' \leftarrow \mathcal{A}^{\text{H}_i, \text{BD}_i}(y_b)$ return $(b' = b)$	for $i = 1, 2$ do $\text{H}_i \leftarrow \text{Fun}[n_i, m_i]$ $(x_1, x_2) \leftarrow \mathcal{A}^{\text{H}_i, \text{BD}_i}$ $y_1 \leftarrow \text{C}^{\text{H}_i}(x_1); y_2 \leftarrow \text{C}^{\text{H}_i}(x_2)$ return $(x_1 \neq x_2 \wedge y_1 = y_2)$

Fig. 1: The one-way, pseudorandomness, and collision resistance games for $\text{C}^{\text{H}_i} \in \text{Fun}[n, m]$.

$\text{H}_2(x) := \text{H}(1|x)$ and \mathcal{F} to consist of two types of functions: those in \mathcal{F}_1 and dependent on values $\text{H}(0|x)$, that is the function table of H_1 , only, and those in \mathcal{F}_2 and dependent on values of $\text{H}(1|x)$, that is the function table of H_2 , only. Thus the adversary in the unrestricted 2-BRO model has less power than in the unrestricted 1-BRO model.

2.2 Cryptographic notions

We recall the basic notions of one-wayness, pseudorandomness, and collision-resistance for a construction $\text{C}^{\text{H}_1, \text{H}_2}$ in the 2-BRO model in Figure 1. We omit the implicit parameters from the subscripts and use C^{H_i} in place of $\text{C}^{\text{H}_1, \text{H}_2}$ to ease notation. These notions can also be defined in the 1-BRO model analogously by removing access to H_2 and BD_2 throughout. The advantage terms are

$$\text{Adv}_{\text{C}^{\text{H}_i}}^{\text{ow}}(\mathcal{A}) := \Pr[\text{OW}_{\text{C}^{\text{H}_i}}^{\mathcal{A}}], \quad \text{Adv}_{\text{C}^{\text{H}_i}}^{\text{prg}}(\mathcal{A}) := 2 \cdot \Pr[\text{PRG}_{\text{C}^{\text{H}_i}}^{\mathcal{A}}] - 1,$$

$$\text{Adv}_{\text{C}^{\text{H}_i}}^{\text{cr}}(\mathcal{A}) := \Pr[\text{CR}_{\text{C}^{\text{H}_i}}^{\mathcal{A}}].$$

All probabilities in this model are also taken over random choices of H_i . Informally $\text{C}^{\text{H}_1, \text{H}_2}$ is OW, PRG, or CR if the advantage of any adversary \mathcal{A} querying its oracles, such that the total length of the received responses remains “reasonable”, is “small”. Note that if one only considers backdoor functions with 1-bit output lengths, the total length of the oracle responses directly translates to the number of queries made by \mathcal{A} . We denote by $\text{Q}(\mathcal{A})$ the number of oracle queries made by an adversary \mathcal{A} to H_i and BD_i . Weak security in each case means that the corresponding advantage is less than 1 and not overwhelming.

We define variants of the above games which will be helpful in our analyses. For a function $\text{H} \in \text{Fun}[n, m]$, define $\text{Im}(\text{H}) := \text{H}(\{0, 1\}^n)$ and $\overline{\text{Im}}(\text{H}) := \{0, 1\}^m \setminus \text{Im}(\text{H})$. The *random preimage-resistance* (rPre) game is defined similarly to everywhere preimage-resistance (ePre) [35] except that a random co-domain point (as opposed to any such point) must be inverted. This definition differs from one-way security in two aspects: the distribution of $\text{H}(x)$ for a uniform x might not be uniform. Furthermore, some points in the co-domain might

Game rPre_C^A	Game oPRG_C^A	Game IU_C^A
for $i = 1, 2$ do $H_i \leftarrow \text{Fun}[n_i, m_i]$ $y \leftarrow \{0, 1\}^m$ $x' \leftarrow \mathcal{A}^{H_i, \text{BD}_i}(y)$ if $y \in \overline{\text{img}}(C^{H_i})$ \quad return $(x' = \perp)$ return $(C^{H_i}(x') = y)$	for $i = 1, 2$ do $H_i \leftarrow \text{Fun}[n_i, m_i]$ $y \leftarrow \{0, 1\}^m$ $b' \leftarrow \mathcal{A}^{H_i, \text{BD}_i}(y)$ return $(b' = (y \in \text{img}(C^{H_i})))$	for $i = 1, 2$ do $H_i \leftarrow \text{Fun}[n_i, m_i]$ $y_1 \leftarrow \text{img}(C^{H_i}); b \leftarrow \{0, 1\}$ $x \leftarrow \{0, 1\}^n; y_0 \leftarrow C^{H_i}(x)$ $b' \leftarrow \mathcal{A}^{H_i, \text{BD}_i}(y_b)$ return $(b' = b)$

Fig. 2: The random preimage resistance (rPre), oblivious PRG (oPRG), and image uniformity (IU) games for $C^{H_i} \in \text{Fun}[n, m]$.

not have any preimages. We also define a decisional variant, called *oblivious PRG* (oPRG), where the adversary has to decide if a random co-domain point has a preimage. We formalize these games in Figure 2. The advantage terms are defined as:

$$\text{Adv}_{C^{H_i}}^{\text{rpre}}(\mathcal{A}) := \Pr[\text{rPre}_{C^{H_i}}^A] \quad \text{Adv}_{C^{H_i}}^{\text{oprpg}}(\mathcal{A}) := \Pr[\text{oPRG}_{C^{H_i}}^A]$$

Weak analogues of the above security notions (for example weak rPre or weak oPRG) are defined by requiring the advantage to be bounded away from 1 (i.e., not to be overwhelming). These definitions can be formalized in the asymptotic language, but we use concrete parameters here.

We state two lemmas that relate OW and rPre, resp. PRG and oPRG: for functions that have *uniform images*, as defined below, we show that OW security is implied by rPre security and PRG security is implied by oPRG security.

IMAGE UNIFORMITY. Let $C^{H_i} \in \text{Fun}[n, m]$ be a construction in the 2-BRO model. In the image uniformity game IU defined in Figure 2, an adversary, given access to all backdoor oracles, must decide whether a given value is chosen uniformly at random from the image of C^{H_i} or computed as the image of a value x chosen uniformly at random from the domain. The advantage term is

$$\text{Adv}_{C^{H_i}}^{\text{iu}}(\mathcal{A}) := 2 \cdot \Pr[\text{IU}_{C^{H_i}}^A] - 1,$$

where the probability is taken over random choices of H_i .

The following lemma upper bounds the advantage of adversaries playing the image uniformity game for combiners with different stretch values. We denote by \mathcal{U}_S the uniform distribution over a set S . We also let \mathcal{U}_f^p denote the distribution defined by $\mathcal{U}_f^p(x) = |f^{-1}(x)|/2^n$, where $f \in \text{Fun}[n, m]$ is a uniform function. We refer the readers to [4, Appendix A] for proofs.

Lemma 1 (Combiner image uniformity). *Let $C_t^{H_1, H_2} : \{0, 1\}^n \rightarrow \{0, 1\}^m$ be a combiner for $t \in \{!, \circ\}$. Let $H : \{0, 1\}^n \rightarrow \{0, 1\}^m$ be a hash function. Then*

$$\text{Adv}_{C_t^{H_i}}^{\text{iu}}(\mathcal{A}) \leq \mathbb{E}_H [\Delta_{\text{TV}}(\mathcal{U}_{\text{img}(H)}, \mathcal{U}_H^p)] + 2 \cdot p_t,$$

where $p_1 = 0$ and $p_0 \leq 2^{2n_1 - m_1}$ is the probability that $H_1 : \{0, 1\}^{n_1} \rightarrow \{0, 1\}^{m_1}$ is not injective (i.e., it has at least one collision). Let $2^n = C \cdot 2^{m \cdot \gamma}$ for constants C and γ . Then the above statistical distance is negligible for $\gamma > 1$ and $0 < \gamma < 1$ when $C = 1$, while for $\gamma = 1$ and $C \leq 1$ it is less than $e^{-C} \cdot (C / (1 - e^{-C}) - 1)$ plus negligible terms.

Now we can relate our notions of rPre and oPRG with their classical variants, i.e., one-way and PRG security. Proofs of both lemmas are included in the full version [4, Appendix B].

Lemma 2 (rPre + IU \implies OW). *Let $C^{H_i} \in \text{Fun}[n, m]$ be a construction in the 2-BRO model. Then for any adversary \mathcal{A} against the one-way security of C^{H_i} , there is an adversary \mathcal{B} against the image uniformity and an adversary \mathcal{C} against the rPre security of C^{H_i} , all in the 2-BRO model and using identical backdoor functionalities, such that*

$$\text{Adv}_{C^{H_i}}^{\text{ow}}(\mathcal{A}) \leq \text{Adv}_{C^{H_i}}^{\text{iu}}(\mathcal{B}) + \frac{1}{\alpha} \cdot \text{Adv}_{C^{H_i}}^{\text{rpre}}(\mathcal{C}) - \frac{1 - \alpha}{\alpha},$$

where $\alpha := \Pr[y \in \text{Img}(C^{H_i})]$ over a random choice of $y \in \{0, 1\}^m$ and H_i .

An analogous result also holds for oPRG security.

Lemma 3 (oPRG + IU \implies PRG). *Let $C^{H_i} \in \text{Fun}[n, m]$ be a construction in the 2-BRO model which is expanding with $m - n \geq 0.53$. Then for any adversary \mathcal{A} against the PRG security of C^{H_i} , there is an adversary \mathcal{B} against the image uniformity and an adversary \mathcal{C} against the oPRG security of C^{H_i} , both in the 2-BRO model and using identical backdoor functionalities, such that*

$$\text{Adv}_{C^{H_i}}^{\text{prg}}(\mathcal{A}) \leq \text{Adv}_{C^{H_i}}^{\text{iu}}(\mathcal{B}) + \frac{1 - \alpha}{\alpha} \cdot \text{Adv}_{C^{H_i}}^{\text{oprpg}}(\mathcal{C}) - (1 - \alpha),$$

where $\alpha := \Pr[y \in \text{Img}(C^{H_i})]$ over a random choice of $y \in \{0, 1\}^m$ and H_i .

3 Black-Box Combiners

A standard way to build a good hash function from a number of possibly “faulty” hash functions is to combine them [25]. For instance, given k hash functions H_1, \dots, H_k , the classical concatenation combiner is guaranteed to be collision resistant as long as one out of the k hash functions is collision resistant. More formally, a black-box collision-resistance combiner \mathcal{C} is a pair of oracle circuits $(C^{H_i}, R^{\mathcal{A}})$ where C^{H_i} is the construction and $R^{\mathcal{A}}$ is a reduction that given as oracle any procedure \mathcal{A} that finds a collision for C^{H_i} , returns collisions for *all* of the underlying H_i ’s. We are interested in a setting where *none* of the available hash functions is good. Under this assumption, however, a secure hash function must be built from scratch, implying that the source of cryptographic hardness must lie elsewhere. As we discussed above, this question has been studied in the RO model.

We briefly explore the difficulty in the standard model here. We consider a variant of this problem where the hash functions are weak due to the existence of *backdoors*. A generation algorithm Gen outputs keys (hk, bk) , where hk is used for hashing and bk enables an unspecified backdoor capability (such as finding preimages or collisions). Our hardness assumption is that the hash function with key hk is collision resistant without access to bk . However, when bk is available, no security is assumed. In this setting, the definition of a combiner can be simplified: instead of requiring the existence of a reduction R^A as above, we can proceed in the standard way and require that the advantage of any adversary $\mathcal{A}(S)$ that gets any subset $S \subset \{bk_1, \dots, bk_k\}$ of the backdoors of size $|S| \leq k - 1$ to be small.⁵ Let us call a combiner secure against any set of at most t backdoors a $\binom{k}{k-t}$ -combiner.

It is trivial to see that a $\binom{k}{0}$ -combiner is also a $\binom{k}{1}$ -combiner. It is also easy to see that a black-box combiner is also a $\binom{k}{1}$ -combiner. We are, however, interested in the feasibility of $\binom{k}{0}$ -combiners. In this setting there *is* an assumed source of hardness, namely the collision resistance of hash functions without backdoors. But constructions that have to work with a *provided* set of keys seem hard.⁶ We next give a simple impossibility result that formalizes this intuition under fully black-box constructions.

Theorem 1. *For any positive $k \in \mathbb{N}$, there are no fully black-box constructions of compressing collision-resistant $\binom{k}{0}$ -combiners.*

Proof idea. Let (H, \mathcal{A}) be a pair of oracles such that $H(hk, \cdot)$ implements a random function and $\mathcal{A}(\langle C \rangle, hk_1, \dots, hk_k, bk_1, \dots, bk_k)$ is a break oracle that operates as follows. It interprets $\langle C \rangle$ as the description of a combiner. It then checks that each bk_i indeed enables generating collisions under hk_i . If so, it (inefficiently) finds a random collision for $C^{H(hk_1, \cdot), \dots, H(hk_k, \cdot)}$ and returns it. An efficient reduction R is given oracle access to \mathcal{A} and H as well as a key hk^* (without its backdoor bk^*). It should find a collision for $H(hk^*, \cdot)$ while making a small (below birthday) number of queries to the two oracles \mathcal{A} and H . We show that any such reduction R must have a negligible success probability.

We distinguish between two cases based on whether the reduction R uses the provided break oracle \mathcal{A} or not. Without the use of \mathcal{A} , the reduction would break collision resistance for hk^* on its own, contradicting the collision resistance of hk^* beyond the birthday bound. To use \mathcal{A} the reduction has to provide it with

⁵ The classical setting can be viewed as one where bk 's are *fixed*, which leads to a difficulty when the new definition is used: a combiner (formally speaking) can “detect” which hash functions are the good ones and use them. Since this detection procedure is not considered practical, one instead asks for the existence of a reduction R as discussed above.

⁶ Without this restriction, a trivial construction exists: generate a fresh hash key and “forget” its backdoor. In practice, however, hash keys model sampling of a (unkeyed) hash function from a family. Moreover, it is unclear if the designer of the combiner will securely erase the generated backdoor. Thus, we assume that for any generated key its backdoor is also available.

k keys hk_i and some other keys bk_i that enable finding collisions (since \mathcal{A} checks this). However, none of the provided keys hk_i can be hk^* , since R must also provide some \tilde{bk}^* that enables finding collisions under hk^* , which means that R can directly use \tilde{bk}^* to compute a collision for $\mathsf{H}(hk^*, \cdot)$, once again contradicting the assumed collision resistance of hk^* beyond the birthday bound. Thus, R does not use hk^* . A random oracle $\mathsf{H}(hk^*, \cdot)$, however, is collision resistant even in the presence of random collisions for $\mathsf{H}(hk, \cdot)$ for $hk \neq hk^*$. This means that R , which places a small number of oracle queries, will have a negligible success probability. \square

There is room to circumvent this result by considering non-black-box constructions. Here, we will study hash function combiners in the k -BRO model, where the hash oracles model access to different hk and the backdoor oracles model access to the corresponding bk 's. As mentioned above, this approach has also been adapted in a number of previous works, both from a provable security as well as a cryptanalytic view [20,31,23,26,21]. In this work we will focus on basic security properties of the concatenation (parallel) and cascade (sequential) combiners in the *unrestricted* 2-BRO model.

4 Communication Complexity

The communication cost [38,24] of a two-party deterministic protocol π on inputs (x, y) is the number of bits that are transmitted in a run of the protocol $\pi(x, y)$. We denote this by $\text{CC}(\pi(x, y))$. The worst-case communication complexity of π is $\max_{(x,y)} \text{CC}(\pi(x, y))$. A protocol π computes a task (function) $f : X \times Y \rightarrow Z$ if the last message of $\pi(x, y)$ is $f(x, y)$. The communication complexity of a task f is the minimum communication complexity of any protocol π that computes f . Protocols can also be randomized and thus might err with probability $\Pr[\pi(x, y) \neq f(x, y)]$. Following cryptographic conventions, we denote protocol correctness by $\text{Adv}_\mu^f(\pi)$, where f is a placeholder for the name of the task f .

In the cryptographic setting we are interested in distributional (aka. average-case) communication complexity measured by averaging the communication cost over random choices of inputs and coins. A standard coin-fixing argument shows that in the *distributional* setting any protocol can be derandomized with no change in communication complexity, and thus we can focus on deterministic protocols. For a given distribution μ over the inputs (x, y) , the protocol error and correctness are computed by taking the probability over the choice of (x, y) . We define the distributional communication cost of a deterministic protocol π as

$$D_\mu(\pi) := \mathbb{E}_{(x,y) \sim \mu} [\text{CC}(\pi(x, y))] .$$

The distributional communication complexity of a task f with error ε is

$$D_\mu^\varepsilon(f) := \min_{\pi} D_\mu(\pi) ,$$

where the minimum is taken over all deterministic protocols π which err with probability at most ε . In this work, we need to slightly generalize functional tasks to relational tasks $R(x, y) \subseteq Z$ and define error as $\Pr[\pi(x, y) \notin R(x, y)]$.

Two central problems in communication complexity that have received substantial attention are the *set-disjointness* and the *set-intersection* problems. In set-disjointness two parties, holding sets S and T respectively, compute the binary function $\text{DISJ}(S, T) := (S \cap T = \emptyset)$. In set-intersection, their goal is to compute the relation $\text{INT}(S, T) := S \cap T$; that is, the last message of the protocol should be equal to some element in the intersection. Note that set-disjointness can be seen as a decisional version of set-intersection and is easier. As mentioned, we are interested in average-case lower bounds for these tasks and moreover we focus on *product* distributions, where the sets are chosen independently.

Two main results to this end have been proven.⁷ A classical result of Babai, Frankl, and Simon [1] establishes an $\Omega(\sqrt{N})$ lower bound for set-disjointness where the input sets S and T are independent random subsets of $[N]$ of size exactly \sqrt{N} . This result, however, is restrictive for us as it roughly translates to regular functions in the cryptographic setting. Moreover, its proof uses intricate combinatorial arguments, which are somewhat hard to work with.

A second result considers the following distribution. Each element $x \in [N]$ is thrown into S independently with probability p . (And similarly for T with probability q .) We can view S as a N -bit string X where its i -th bit x_i is 1 iff $i \in S$. Thus the distribution can be viewed as N i.i.d. Bernoulli random variables $x_i \sim \text{Ber}(p)$ where $p := \Pr[x_i = 1]$. Thus the elements of the sets form a binomial distribution, and accordingly we write $S \sim \text{Bin}(N, p)$ and $T \sim \text{Bin}(N, q)$. We define $\mu(p, q)$ as the *product* of these distributions. When $p = q = 1/2$ we get the product uniform distribution over the subsets of $[N] \times [N]$, but typically we will be looking at much smaller values of p and q of order $1/\sqrt{N}$.

Using information-theoretic techniques [2], the following lower bound can be established.

Theorem 2 (Set-Disjointness Lower Bound). *Let $N \in \mathbb{N}$ and assume $p, q \in (0, 1/2]$ with $p \leq q$ and $pq = 1/(\delta N)$ for some $\delta > 1$. Let $\mu(p, q)$ be the product binomial distribution over subsets $S, T \subseteq [N]$. Assume $\varepsilon < \frac{(\delta-1)p_0}{(4+\delta)}$ and let $p_0 := \Pr[\text{DISJ}(S, T) = 0]$. Then*

$$D_{\mu(p, q)}^{\varepsilon}(\text{DISJ}) \geq \frac{Np}{8} \cdot ((\delta - 1)p_0 - (4 + \delta)\varepsilon)^2.$$

We have included a detailed proof of the above theorem in the full version [4, Appendix C], which follows those in [32, 19]. Our proof generalizes the original result, which was only claimed for $p = q = 1/\sqrt{N}$.⁸ Roughly speaking, the

⁷ We note that most of the work on distributional communication complexity is driven by Yao's min-max lemma, which lower bounds worst-case communication complexity using distributional communication complexity for *some* (often non-uniform) distribution.

⁸ In the full version of this paper [4, Appendix C.3] we give a new refined proof that extends the theorem to $\delta \geq 0.8$.

proof proceeds along the following lines. We can lower bound the communication complexity of any protocol by the total information leaked by its transcripts about each coordinate (x_i, y_i) . The latter can be lower bounded based on the statistical distance in protocol transcripts when $x_i = 1 \wedge y_i = 0$ and $x_i = 0 \wedge y_i = 1$. This step uses a number of information-theoretic inequalities, which we include with proofs in the full version. Finally, we show that a highly correct protocol can be used as a distinguisher with constant advantage: When $x_i = 0 \wedge y_i = 0$, for a constant fraction of the inputs the sets will be disjoint. However, when $x_i = 1 \wedge y_i = 1$ they necessarily intersect, but this condition happens for a constant fraction of the inputs. We get a \sqrt{N} lower bound by averaging over the i 's.

In this section we also prove a communication complexity lower bound for the set-intersection problem over Bernoulli sets for which set-disjointness can be *easy*. Although the overall proof structure will be similar to that in [32,19], we will differ in a number of places. First, as above we leave the Bernoulli parameters free so as to be able to compute a feasible region where the lower bound will be non-trivial. We also use the fact that a candidate element can be checked to belong to the intersection (whereas a decision bit for disjointness cannot be checked for correctness). This ensures that the protocol error is one-sided, and allows us to remove the requirement of ε being sufficiently small. Finally, we will bound the probability that the protocol outputs a *random element* in the intersection. This leads to a distinguisher that succeeds with smaller advantage, but overall will lead to a non-trivial bound. We state and prove the formal result next.

Theorem 3 (Set-Intersection Lower Bound). *Let $N \in \mathbb{N}$ and assume $p, q \in (0, 1/2]$ with $p \leq q$. Let $\mu(p, q)$ be the product binomial distribution over subsets $S, T \subseteq [N]$. Let ε be the protocol error and set $p_0 := \Pr[\text{DISJ}(S, T) = 0]$. If $\varepsilon \leq p_0$ then*

$$D_{\mu(p,q)}^{\varepsilon}(\text{INT}) \geq \frac{Np}{8} \cdot \left(\frac{p_0 - \varepsilon}{Npq} \right)^2 .$$

For sufficiently large N we have $p_0 = 1 - (1 - pq)^N \approx 1 - e^{-Npq}$. If $pq \gg 1/N$ we get that $p_0 \approx 1$ (the sets intersect with overwhelming probability) and for the theorem we would need that $\varepsilon \leq 1$.

Let us first give some preliminaries and state two lemmas that are used in the proof of Theorem 3. For random variables X and Y , their statistical distance (aka. total variance) is denoted by $\Delta_{\text{TV}}(X, Y)$, their mutual information is denoted by $I(X; Y)$, and their Hellinger distance is denoted by $\Delta_{\text{Hel}}(X, Y)$:

$$\Delta_{\text{Hel}}(X, Y) := \sqrt{1 - \sum_{z \in D} \sqrt{\Pr[X = z] \Pr[Y = z]}} .$$

Statistical and Hellinger distance are related (cf. proofs in [4, Appendix C.1]) via:

$$\Delta_{\text{Hel}}^2(X, Y) \leq \Delta_{\text{TV}}(X, Y) \leq \sqrt{2} \cdot \Delta_{\text{Hel}}(X, Y) .$$

Below, Lemma 4, proven in the full version [4, Appendix C.1], relates the mutual information of two random variables with their Hellinger distance.

Lemma 4 (Information to Hellinger). *Let X and Y be random variables and $Y_x := Y|X = x$, i.e., Y conditioned on $X = x$. Then*

$$\mathbb{E}_{x \in X} [\Delta_{\text{Hel}}^2(Y, Y_x)] \leq I(X; Y) .$$

Next we state the cut-and-paste lemma from communication complexity. A proof is included in [4, Appendix C.2].

Lemma 5 (Cut-and-Paste). *Let $\Pi(X, Y)$ denote a random variable for the transcripts of a deterministic protocol on input bit strings (X, Y) such that the corresponding sets S and T are drawn from μ , i.e., $S, T \sim \mu$. Let $a, b \in \{0, 1\}$ and define $\Pi_{a,b}^i(X, Y) := \Pi(X, Y) \mid x_i = a \wedge y_i = b$. Then for each i , it holds that*

$$\Delta_{\text{Hel}}^2(\Pi_{0,0}^i, \Pi_{1,1}^i) = \Delta_{\text{Hel}}^2(\Pi_{0,1}^i, \Pi_{1,0}^i) .$$

Now we can prove the claimed lower bound on the communication complexity of set-intersection.

Proof of Theorem 3. Let π be a deterministic protocol with error at most ε , i.e.,

$$\Pr_{(S,T) \sim \mu} [\pi(X, Y) \in \text{INT}(S, T)] \geq 1 - \varepsilon ,$$

where X and Y are bit string representations of S and T as explained above. Let $\Pi(X, Y)$ denote a random variable for the transcripts of protocol π on inputs (X, Y) with corresponding sets $(S, T) \sim \mu$. We write $X = (x_1, \dots, x_N)$ and $Y = (y_1, \dots, y_N)$ where $x_i, y_i \in \{0, 1\}$. For random variables A and B , let $\text{supp}(A)$ denote the support of A (i.e., the set of values that have a non-zero probability of happening), and $H(A)$ denote the Shannon entropy. We have

$$\begin{aligned} D_{\mu(p,q)}^\varepsilon(\text{INT}) &\geq \log |\text{supp}(\Pi(X, Y))| \\ &\geq H(\Pi(X, Y)) = H(\Pi(X, Y)) + H(X, Y) - H(X, Y, \Pi(X, Y)) \\ &= I(X, Y; \Pi) = I(x_1, \dots, x_N, y_1, \dots, y_N; \Pi) \geq \sum_{i=1}^N I(x_i, y_i; \Pi) , \end{aligned}$$

where the last inequality holds due to the independence of $x_1, \dots, x_N, y_1, \dots, y_N$ (cf. [4, Appendix C.1]). Let $\Pi_{a,b}^i$ be Π conditioned on the i -th coordinates of X and Y being fixed to a and b respectively:

$$\Pi_{a,b}^i(X, Y) := \Pi(X, Y) \mid x_i = a \wedge y_i = b .$$

By Lemma 4 we know

$$I(x_i, y_i; \Pi) \geq \mathbb{E}_{(a,b)} [\Delta_{\text{Hel}}^2(\Pi_{a,b}^i, \Pi)] ,$$

where $(a, b) \sim \text{Ber}(p) \times \text{Ber}(q)$ and Δ_{Hel} is the Hellinger distance.

Since $q \geq p$ we have that $q(1-p) \geq p(1-q)$ and since $q \leq 1/2$, we also have that $p(1-q) \geq p/2$. Thus

$$\begin{aligned}
I(x_i, y_i; \Pi) &\geq p(1-q) \cdot \Delta_{\text{Hel}}^2(\Pi_{1,0}^i, \Pi) + q(1-p) \cdot \Delta_{\text{Hel}}^2(\Pi_{0,1}^i, \Pi) \\
&\geq p(1-q) \cdot (\Delta_{\text{Hel}}^2(\Pi_{1,0}^i, \Pi) + \Delta_{\text{Hel}}^2(\Pi_{0,1}^i, \Pi)) \\
&\geq p/2 \cdot (\Delta_{\text{Hel}}^2(\Pi_{1,0}^i, \Pi) + \Delta_{\text{Hel}}^2(\Pi_{0,1}^i, \Pi)) \\
&\geq p/4 \cdot (\Delta_{\text{Hel}}(\Pi_{1,0}^i, \Pi) + \Delta_{\text{Hel}}(\Pi_{0,1}^i, \Pi))^2 \\
&\geq p/4 \cdot \Delta_{\text{Hel}}^2(\Pi_{1,0}^i, \Pi_{0,1}^i) .
\end{aligned}$$

The last inequality is by the triangle inequality for the metric Δ_{Hel} , and the penultimate inequality uses $x^2 + y^2 \geq (x+y)^2/2$. Hence,

$$\begin{aligned}
D_{\mu(p,q)}^\varepsilon(\text{INT}) &\geq N \cdot \mathbb{E}_i[I(x_i, y_i; \Pi)] \\
&\geq Np/4 \cdot \mathbb{E}_i[\Delta_{\text{Hel}}^2(\Pi_{1,0}^i, \Pi_{0,1}^i)] \\
&= Np/4 \cdot \mathbb{E}_i[\Delta_{\text{Hel}}^2(\Pi_{0,0}^i, \Pi_{1,1}^i)] \\
&\geq Np/8 \cdot \mathbb{E}_i[\Delta_{\text{TV}}^2(\Pi_{0,0}^i, \Pi_{1,1}^i)] \\
&\geq Np/8 \cdot (\mathbb{E}_i[\Delta_{\text{TV}}(\Pi_{0,0}^i, \Pi_{1,1}^i)])^2 ,
\end{aligned}$$

where the third inequality uses the cut-and-paste lemma of communication complexity (Lemma 5) which states that $\Delta_{\text{Hel}}^2(\Pi_{1,0}^i, \Pi_{0,1}^i) = \Delta_{\text{Hel}}^2(\Pi_{0,0}^i, \Pi_{1,1}^i)$ for any deterministic protocol π . The penultimate inequality uses $\Delta_{\text{TV}}(A, B) \leq \sqrt{2}\Delta_{\text{Hel}}(A, B)$, which implies $\Delta_{\text{Hel}}^2(\Pi_{0,0}^i, \Pi_{1,1}^i) \geq 1/2\Delta_{\text{TV}}^2(\Pi_{0,0}^i, \Pi_{1,1}^i)$, and the last inequality is by Jensen. Thus it remains to lower bound $\Delta_{\text{TV}}(\Pi_{0,0}^i, \Pi_{1,1}^i)$.

For every i we have

$$\Pr[\Pi_{0,0}^i(X, Y) = i] = 0 .$$

This is because we have conditioned on $x_i = y_i = 0$ and the two parties can check whether or not i belongs to their sets.

Now we look at $x_i = y_i = 1$. We show that the protocol over a *random* choice of i should output i with the expected probability, that is, $1/|S \cap T|$. Note that the expected size of the intersection is

$$\mathbb{E}[|S \cap T|] = \mathbb{E}\left[\sum_{i=1}^N x_i y_i\right] = \sum_{i=1}^N \mathbb{E}[x_i y_i] = Npq ,$$

where we have used the linearity of expectation and independence of x_i and y_i . We proceed as follows.

$$\begin{aligned}
\mathbb{E}_i[\Pr[\Pi_{1,1}^i(X, Y) = i]] &= \\
&= \frac{1}{N} \sum \Pr[\Pi(X, Y) = i | x_i = y_i = 1]
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{Npq} \sum \Pr[\Pi(X, Y) = i \wedge x_i = y_i = 1] \\
&= \frac{1}{Npq} \sum_i \sum_{(x, y): x_i = y_i = 1 \wedge \pi(x, y) = i} \Pr[(X, Y) = (x, y)] \\
&= \frac{1}{Npq} \sum_{(x, y): \pi(x, y) \text{ correct and } x \cap y \neq \emptyset} \Pr[(X, Y) = (x, y)] \\
&= \frac{1}{Npq} \left(\sum_{(x, y)} \Pr[(x, y)] - \sum_{x \cap y = \emptyset} \Pr[(x, y)] - \sum_{\pi(x, y) \text{ fails}} \Pr[(x, y)] \right) \\
&\geq \frac{1 - \Pr[\text{DISJ}(S, T) = 1] - \varepsilon}{Npq} = \frac{p_0 - \varepsilon}{Npq}.
\end{aligned}$$

Thus we get that

$$\mathbb{E}_i[\Delta_{\text{TV}}(\Pi_{0,0}^i, \Pi_{1,1}^i)] \geq (p_0 - \varepsilon)/(Npq),$$

and overall we obtain

$$D_{\mu(p,q)}^\varepsilon(\text{INT}) \geq \frac{Np}{8} \cdot \left(\frac{p_0 - \varepsilon}{Npq} \right)^2,$$

as required. □

Letting $p = 1/N^\alpha$ and $q = 1/N^\beta$ with $\alpha \geq \beta$ (since we assumed $p \leq q$), for a non-trivial lower bound—that is an exponentially large right-hand side in the displayed equation above—we would need to have that $\alpha + 2\beta > 1$. We also require that $1 - \alpha - \beta > 0$ so that the expected intersection size Npq is exponentially large, in which case $p_0 \approx 1$ and set-disjointness is *easy*. These inequalities lead to the feasibility region shown in Figure 3. We have included the symmetric region for $\alpha \leq \beta$.

In this work, we will rely on set-disjointness and set-intersection problems, as well as the following *multi-set* extensions of them. These problems are additionally parameterized by the number of sets. Here Alice holds M_1 sets $S_i \sim \text{Bin}(N, p)$ for $i \in [M_1]$ and Bob holds M_2 sets $T_j \sim \text{Bin}(N, q)$ for $j \in [M_2]$. Their goal is to solve the following problems.

1. Find (i, x) such that $x \in S_i \cap T_i$, or return \perp if all the intersections are empty. We call this the (M_1, M_2) -INT problem, a natural multi-instance version of INT. A decisional variant would ask for an index i and a decision bit indicating if $S_i \cap T_i = \emptyset$. When $M_1 = M_2 = 1$, these problems are the usual INT and DISJ problems.

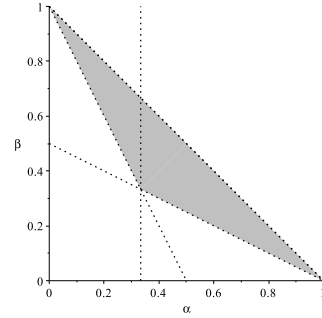


Fig. 3: Region where set intersection is hard with $p = 1/N^\alpha$ and $q = 1/N^\beta$.

2. Find (i, j, x, x') with $x \neq x'$ such that $x, x' \in S_i \cap T_j$, or return \perp if no such tuple exists. We call this the (M_1, M_2) -2INT problem. When $M_1 = M_2 = 1$ this problem is at least as hard as the INT problem since finding two distinct elements in the intersection is harder than finding one.

REMARK. Intuitively, the INT problem is a harder task than (M_1, M_2) -2INT. One can solve the (M_1, M_2) -2INT problem using a protocol for INT as follows. Alice chooses a random point x in one of its sets S_i and sends it to Bob. Bob will then search through his sets to find a set T_j such that $x \in T_j$. With high probability such a set exists if the number of sets and/or the probability parameters are large enough. Alice and Bob will then run the protocol for INT on sets S_i and T_j to find an $x' \in S_i \cap T_j$. This element will be different from x with good probability (again under appropriate choices of parameters). Indeed, this is simply the communication complexity way of saying “collision-resistance implies one-wayness.” However, we are interested in a reduction in the converse direction (as we already have lower bounds for INT). This seems hard as from a cryptographic point of view, as a classical impossibility by Simon [36] shows that collision resistance cannot be based on one-way functions (or even permutations) in a black-box way. Despite this, it is conceivable that direct information-theoretic analyses (similar to those for set-disjointness and set-intersection) can lead to non-trivial lower bounds. We leave proving hardness for this “collision resistance” analogue of set-intersection as an interesting open problem for future work.⁹

5 The Concatenation Combiner

In this section we study the security of the concatenation combiner

$$\mathbf{C}^{\mathbf{H}_1, \mathbf{H}_2}(x) := \mathbf{H}_1(x) \parallel \mathbf{H}_2(x)$$

in the 2-BRO model, where $\mathbf{H}_1 \in \text{Fun}[n, n + s_1]$ and $\mathbf{H}_2 \in \text{Fun}[n, n + s_2]$. We will prove one-way security, pseudorandomness, and collision resistance for this construction. Our results will rely on the hardness of set-intersection and set-disjointness for the first two properties, and the presumed hardness of finding two elements in the intersection given multiple instances.

5.1 One-way security

In the full version of this paper [4, Appendix D] we show that when \mathbf{H}_1 or \mathbf{H}_2 is (approximately) length preserving or somewhat expanding the concatenation combiner is not (strongly) one-way in the 2-BRO model. In both cases preimage

⁹ We note that our setting is different to direct sum/product theorems where the focus is on hardness amplification. One shows, for example, that computing n independent copies of a function requires n times the communication for one copy for product distributions [3].

sets will be only polynomially large and can be communicated. Accordingly, only when both hash functions are (somewhat) compressing we can achieve one-way security.

To this end, we first give a direct reduction from random preimage resistance (rPre, as defined in Figure 2) to set-intersection. By Lemma 2 we know that any (weak) rPre-secure function is also a (weak) OW-secure function. In particular, for the highly compressing setting where $s_1, s_2 \leq -n/2 - 4$ we show *strong* one-way security. For settings where the parameters only enable weak security according to the set-intersection theorem, we can apply hardness amplification [18] to get a strongly one-way function.

In our reductions to communication complexity protocols throughout the paper, we make the following simplifying assumptions. 1) The adversary is deterministic; 2) It does not query H_i at all and instead computes hash values via the BD_i oracles; 3) It queries BD_i with functions that have 1-bit outputs; and 4) It starts with a query to BD_1 .

We are now ready to prove our first cryptographic hardness result.

Theorem 4. *Let $H_1 \in \text{Fun}[n, n + s_1]$ and $H_2 \in \text{Fun}[n, n + s_2]$ and $C^{H_1, H_2}(x) := H_1(x)|H_2(x)$. Then for any adversary \mathcal{A} against the rPre security of C^{H_1, H_2} in the 2-BRO model there is a 2-party protocol π against set-intersection with $\mu := \mu(p, q)$ where $p := 1/2^{n+s_1}$ and $q := 1/2^{n+s_2}$ and such that*

$$\text{Adv}_{C^{H_1, H_2}}^{\text{rpre}}(\mathcal{A}) \leq \text{Adv}_{\mu}^{\text{int}}(\pi) \quad \text{and} \quad D_{\mu}(\pi) \leq Q(\mathcal{A}) + 3n + s_1 + s_2 .$$

Proof. Let \mathcal{A} be an adversary against the rPre security of C^{H_1, H_2} in the 2-BRO model for H_1 and H_2 . Adversary \mathcal{A} is given a random point $y := y_1|y_2 \in \{0, 1\}^{2n+s_1+s_2}$ and needs to either find an x such that $H_1(x)|H_2(x) = y_1|y_2$ or say that no such x exists. Let

$$S_1 := H_1^{-1}(y_1) \quad \text{and} \quad S_2 := H_2^{-1}(y_2) .$$

Hence \mathcal{A} outputs an $x \in S_1 \cap S_2$ as long as $S_1 \cap S_2 \neq \emptyset$. We note that these sets are Bernoulli. Indeed, for each x we have that $\Pr[x \in S_1] = 1/2^{n+s_1}$ and $\Pr[x \in S_2] = 1/2^{n+s_2}$, and these events are independent for different values of x .

We use \mathcal{A} to build a 2-party protocol for set-intersection over a product distribution $\mu := \mu(p, q)$ with $p := 1/2^{n+s_1}$ and $q := 1/2^{n+s_2}$ as follows. Alice holds a set $S_1 \subseteq \{0, 1\}^n$ and Bob holds a set $S_2 \subseteq \{0, 1\}^n$ distributed according to μ . Alice (resp., Bob) samples hash function H_1 (resp., H_2) as follows. Alice picks a random $y_1 \in \{0, 1\}^{n+s_1}$ and Bob picks a random $y_2 \in \{0, 1\}^{n+s_2}$. Alice defines H_1 to map all points in S_1 to y_1 . She maps $x \in \{0, 1\}^n \setminus S_1$ to random points in $\{0, 1\}^{n+s_1} \setminus \{y_1\}$. Similarly Bob defines H_2 to map all points $x \in S_2$ to y_2 and $x \in \{0, 1\}^n \setminus S_2$ to random points in $\{0, 1\}^{n+s_2} \setminus \{y_2\}$. As a result, Alice knows the full function table of H_1 and similarly Bob knows the full function table of H_2 .

Alice and Bob now run two copies of \mathcal{A} in tandem as follows, where the state values st_A and st_B are initially set to $y_1|y_2$ (with only $2n + s_1 + s_2$ bits of communication).

Alice: It resumes/starts $\mathcal{A}(st_A)$. It terminates if it receives a final guess x from Bob. It answers all pending BD_2 queries—there are none to start with—using the values just received from Bob. It answers all BD_1 queries using the function table of H_1 until \mathcal{A} queries BD_2 or terminates. If \mathcal{A} terminates with a final guess x , it forwards x to Bob and terminates. Else it saves the current state st_A of \mathcal{A} locally and forwards all BD_1 answers that it has provided to \mathcal{A} since the last resumption to Bob. It hands the execution over to Bob.

Bob: It resumes $\mathcal{A}(st_B)$. It terminates if it receives a final guess x from Alice. It answers all pending BD_1 queries using the values received from Alice. It answers all BD_2 queries using the function table of H_2 until \mathcal{A} queries BD_1 or terminates. If \mathcal{A} terminates with a final guess x , it forwards x to Alice and terminates. Else it saves the current state st_B of \mathcal{A} locally and forwards all BD_2 answers that it has provided to \mathcal{A} since the last resumption to Alice. It hands the execution over to Alice.

We claim that Alice and Bob run \mathcal{A} in an environment that is identical to the rPre game in the 2-BRO model. The hash functions H_1 and H_2 sampled by Alice and Bob are uniformly distributed. To see this note that for any (x, y) the probability that $H_1(x) = y$ is $1/|\{0, 1\}^{n+s_1}|$ (and similarly for H_2). Furthermore, this event is independent of the hash values that are set for all other values $x' \neq x$. Thus, Alice and Bob faithfully run \mathcal{A} in the environment that it expects by answering its backdoor queries using their knowledge of the full tables of the two functions.

Whenever \mathcal{A} succeeds in breaking the rPre security of C^{H_1, H_2} , the protocol above computes an $x \in S_1 \cap S_2$ or says that no such x exists. In either case, the protocol solves the set-intersection problem. Thus the correctness of this protocol is at least the advantage of the adversary \mathcal{A} .

This execution of \mathcal{A} by Alice and Bob ensures that oracle *queries* do not affect the communication cost of Alice and Bob. It is only their answers (plus the final x) that affects the communication cost, since the queried functions f are locally computed and only their answers are communicated. If \mathcal{A} makes $Q(\mathcal{A})$ queries to BD_1 and BD_2 in total and each query has a 1-bit output, the total communication complexity of the protocol is $Q(\mathcal{A})$ plus those bits needed to communicate y_1 and y_2 and the final guess x . \square

We now check that the parameters for hash functions can be set such that their concatenation is a one-way function.

Corollary 1. *For $H_1, H_2 \in \text{Fun}[n, (1 - \epsilon)n/2]$ with $0 < \epsilon < 1/3$ the concatenation combiner is a strongly one-way compressing function in $\text{Fun}[n, (1 - \epsilon)n]$.*

Proof. The feasible region in Figure 3 for $\alpha = \beta$ consists of $1/3 < \alpha < 1/2$. In our setting $\alpha = \beta = (1 - \epsilon)/2$, which means concatenation is strongly rPre secure when $0 < \epsilon < 1/3$. Since the combined function is compressing (where $\gamma = 1/(1 - \epsilon) > 1$), the image-uniformity bound is negligible and also $\Pr[y \in \text{img}(C^{H_i})]$ in Lemma 2 is overwhelming. Using these bounds and Lemma 2 we get that strong rPre security implies strong OW security. \square

We conjecture that concatenation is strongly one-way even for $1/3 \leq \epsilon < 1$. The intuition is that in the one-way game a point is “planted” in a large intersection, which seems hard to discover without essentially communicating the entire intersection. Tighter lower bounds for set-intersection can be used to establish this.

5.2 PRG security

We now consider the PRG security of the concatenation combiner. Our reduction in Theorem 4 from rPre to set-intersection can be easily adapted to the decisional setting. That is, we can show that a decisional variant of rPre can be reduced to the set-disjointness problem. The decisional variant of rPre asks the adversary to decide whether or not a random co-domain point $y_1|y_2$ has a preimage. This is exactly the oblivious PRG (oPRG) notion that we defined in Section 2. We get the following result.

Theorem 5. *Let $H_1 \in \text{Fun}[n, n + s_1]$ and $H_2 \in \text{Fun}[n, n + s_2]$ and $C^{H_1, H_2}(x) := H_1(x)|H_2(x)$. Then for any adversary \mathcal{A} against the oblivious PRG security of C^{H_1, H_2} in the 2-BRO model there is a 2-party protocol π against set-disjointness with $\mu := \mu(p, q)$ where $p := 1/2^{n+s_1}$ and $q := 1/2^{n+s_2}$ and such that*

$$\text{Adv}_{C^{H_1, H_2}}^{\text{oPRG}}(\mathcal{A}) \leq \text{Adv}_{\mu}^{\text{disj}}(\pi) \quad \text{and} \quad D_{\mu}(\pi) \leq Q(\mathcal{A}) + 2n + s_1 + s_2 + 1 .$$

We next check if concrete parameters can be set to obtain an expanding PRG.

Corollary 2. *For $s_1 = -n/2 + 1$ and $s_2 = -n/2$, the concatenation combiner gives a weak PRG in $\text{Fun}[n, n + 1]$.*

Proof. The theorem gives a reduction to set-disjointness with parameters $p = 1/2^{n/2+1}$ and $q = 1/2^{n/2}$. For large n we get, $\delta = 2$, $p_0 = 1 - e^{-1/2}$ and $(\delta - 1)p_0/(4 + \delta) < 0.0656$, which means we can set $\epsilon = 0.065$. By set-disjointness lower bound, this means any adversary with advantage at least 0.935 must place at least $\mathcal{O}(2^{n/2})$ queries in total to its oracles.

By Lemma 1 we have that $\text{Adv}_{C^{H_i}}^{\text{iu}}(\mathcal{B}) \leq e^{-C} \cdot (C/(1 - e^{-C}) - 1)$. In our case $C = 1/2 < 1$, and the right hand side above is upper bounded by ≤ 0.165 . (We have removed the negligible terms and instead approximated the constants by slightly larger values.)

In Lemma 3 in order to meet the bound $\text{Adv}_{C^{H_i}}^{\text{oPRG}}(\mathcal{C}) < (2 - \alpha - \text{Adv}_{C^{H_i}}^{\text{iu}}(\mathcal{B})) \cdot \alpha/(1 - \alpha)$, we would need $0.935 \leq (2 - \alpha - 0.165) \cdot \alpha/(1 - \alpha)$. After some algebra this gives $\alpha \geq 0.39343$. With $m = n + s$, we need to have $1 - e^{-2^{-s}} \geq 0.39343$, which means $s \leq 1.00018$. Thus we can set $s = 1$ (which also satisfies $s \geq 0.53$ as required in the lemma). \square

We can obtain a strong PRG by amplification. However, we need an amplifier that works on PRGs with (very) *small* stretch. Such a construction is given

by Maurer and Tessaro [29]. In their so-called Concatenate-and-Extract (CaE) construction one sets

$$\text{PRG}(r, x_1, \dots, x_m) := r | \text{Ext}(r, C^{\text{H}_1, \text{H}_2}(x_1)) | \dots | C^{\text{H}_1, \text{H}_2}(x_m) ,$$

where Ext is a sufficiently good randomness extractor, for instance a two-universal hash function. We refer to the original work for concrete parameters. It is safe to assume the extractor is backdoor-free, since it is an information-theoretic object and relatively easy to implement.

5.3 Collision resistance

The classical result of Simon [36] shows that collision-resistance relies on qualitatively stronger assumptions than one-way functions. In the theorem below we prove collision resistance based on the hardness of the multi-instance 2INT problem as defined in Section 4. As discussed in the final remark of that section, we do not expect that a reduction to the INT problem exists.

Theorem 6. *Let $\text{H}_1 \in \text{Fun}[n, n + s_1]$ and $\text{H}_2 \in \text{Fun}[n, n + s_2]$ and $C^{\text{H}_1, \text{H}_2}(x) := \text{H}_1(x) | \text{H}_2(x)$. Then for any adversary \mathcal{A} against the collision resistance of $C^{\text{H}_1, \text{H}_2}$ in the 2-BRO model there is a 2-party protocol π' against multi-instance two-element set-intersection problem over $\mu' := \mu(p', q')$ with $p' := 2n \ln 2 / 2^{n+s_1}$ and $q' := 2n \ln 2 / 2^{n+s_2}$ and where Alice holds $M_1 := 2^{n+s_1}$ sets and Bob holds $M_2 := 2^{n+s_2}$ sets such that*

$$\text{Adv}_{C^{\text{H}_1, \text{H}_2}}^{\text{cr}}(\mathcal{A}) \leq \text{Adv}_{\mu'}^{\text{mi-2int}}(\pi') + 2 \cdot 2^{-n} \quad \text{and} \quad \text{D}_{\mu'}(\pi') \leq \text{Q}(\mathcal{A}) + 4n + s_1 + s_2 .$$

Proof. We follow an overall strategy that is similar to one for the rPre reduction. For each $i \in \{0, 1\}^{n+s_1}$, Alice sets $\text{H}_1^-(i) := S_i$ and for each $j \in \{0, 1\}^{n+s_2}$ Bob sets $\text{H}_2^-(j) := T_j$ and they simulate the two hash functions. However, this leads to a problem: S_i are not necessarily disjoint and furthermore their union does not cover the entire domain $\{0, 1\}^n$. (The same is true for T_j .) Put differently, the distributions of sets formed by hash preimages of co-domain points do not match independently chosen sets from a Bernoulli distribution.

We treat this problem in two step. The first step is a direct reduction to a “partitioned” modification of the multi-instance set-intersection. In this partitioned problem Alice gets sets $S_i := \text{H}_1^-(i)$ for $i \in \{0, 1\}^{n+s_1}$ and a random oracle $\text{H}_1 \in \text{Fun}[n, n + s_1]$. Similarly, Bob gets sets $T_j := \text{H}_2^-(j)$ for $j \in \{0, 1\}^{n+s_2}$ and an independent random oracle $\text{H}_2 \in \text{Fun}[n, n + s_2]$. Their goal is to find a tuple (i, j, x, x') with $x \neq x'$ such that $x, x' \in S_i \cap T_j$. Thus, these sets exactly correspond to hash preimages as needed in the reduction above, and a solution would translate to a collision for the combined hash function.

We then show that hardness of the (standard) multi-instance two-element set-intersection problem implies the hardness of the partitioned problem with an increase in the Bernoulli *parameter*.¹⁰ \square

¹⁰ Another strategy would be to change the *number* of sets involved. But this runs into a problem as this number must match the size of the co-domain of the hash function.

<p style="margin: 0;">Algo. ReDist(S_1, \dots, S_M)</p> <hr style="border: 0; border-top: 1px solid black; margin: 2px 0;"/> <p style="margin: 0;">for $x \in \{0, 1\}^n$ do</p> <p style="margin: 0; padding-left: 20px;">$A_x := \{i \in [M] : x \in S_i\}$</p> <p style="margin: 0; padding-left: 20px;">$i_x \leftarrow A_x$</p> <p style="margin: 0; padding-left: 20px;">for $j \in [M] \wedge j \neq i_x$ do</p> <p style="margin: 0; padding-left: 40px;">$\tilde{S}_j \leftarrow \tilde{S}_j \setminus \{x\}$</p> <p style="margin: 0;">return $(\tilde{S}_1, \dots, \tilde{S}_M)$</p>

Fig. 4: Redistribution of elements to form a partition.

Lemma 6 (Partitioned \implies Independent). *For any two-party protocol π against the partitioned multi-instance set-intersection problem there is a two-party protocol π' against multi-instance set-intersection problem such that*

$$\text{Adv}_{\mu'}^{\text{mi-2int}}(\pi') \geq \text{Adv}_{\mu}^{\text{part-2int}}(\pi) - 2 \cdot 2^{-n} \quad \text{and} \quad \text{D}_{\mu'}(\pi') \leq \text{D}_{\mu}(\pi) .$$

Here $\mu := \mu(p, q)$ is the distribution induced by hash preimages and $\mu' := \mu'(p', q')$ is a product Bernoulli with $p' := 2n \ln 2 \cdot p$ and $q' := 2n \ln 2 \cdot q$.

Proof. To focus on the core ideas, we simplify and let $M_1 = M_2 = M = 2^{n+s}$ and $p = q = 1/2^{n+s}$. Suppose we have sets S_i and T_j for $i = 1, \dots, M$ and $j = 1, \dots, M$ as an instance for the multi-instance intersection. Let $p' = q' = 2n \ln 2 \cdot p$. Then

$$\begin{aligned} \Pr[\exists x \in \{0, 1\}^n \forall i \in [M] : x \notin S_i] &\leq 2^n \Pr[\forall i \in [M] : x \notin S_i] \\ &\leq 2^n (1 - p')^{1/p} \leq 2^n e^{-2n \ln 2} = 2^{-n} . \end{aligned}$$

Thus with these parameters the sets S_i (and similarly T_j) will cover the full domain, that is $\bigcup_{i=1}^M S_i = \{0, 1\}^n$.

Note that with these parameters any two sets S_i and T_j will intersect with overwhelming probability. However, finding an element in the intersection may still be hard; see conjecture below.

Our next step it to redistribute the elements among the sets so that they form partitions. We do this via the algorithm ReDist shown in Figure 4. ReDist iterates through elements x in the domain and leaves x in exactly one of the sets. (By the above covering property such a set always exists.)

This procedure will be applied to S_i (resp., T_j) to produce non-overlapping sets \tilde{S}_i (resp. \tilde{T}_j). Furthermore, we always have that $\tilde{S}_i \subseteq S_i$ and $\tilde{T}_i \subseteq T_i$, since elements are only deleted from the sets and never added to them. Thus $\tilde{S}_i \cap \tilde{T}_j \subseteq S_i \cap T_j$ as well, and this means that any solution with respect to the tweaked sets will also be a valid solution for the original (Bernoulli) sets.

We still need to show that the distribution of the tweaked sets is identical to that given by hash preimages under a random oracle. Let $E_{x,i}$ be the event

that $x \in \tilde{S}_i$. Since the algorithm does not treat any of the i 's in a special way, we claim that $\Pr[E_{x,i}]$ is independent of i . Indeed for any i, j we have

$$\Pr[E_{x,i}] = \Pr[x \in S_i] \Pr[i_x = i | x \in S_i] = \Pr[x \in S_j] \Pr[i_x = j | x \in S_j] = \Pr[E_{x,j}].$$

This is because $\Pr[x \in S_i] = \Pr[x \in S_j]$ and $\Pr[i_x = i | x \in S_i] = \Pr[i_x = j | x \in S_j]$. If we call this common probability e_x , since x is guaranteed belongs to one of the M sets, we have that $\sum_{i \in [M]} e_x = 1$. Thus $e_x = 1/M = \Pr[\mathbf{H}_1(x) = i]$. Note that the algorithm assigns different values of x independently of all other values already assigned, we get that the event $\mathbf{H}_1(x) = i$ is independent for different x .

Finally, solutions with respect to the tweaked sets always exist when $s_1 + s_2 < 0$. This is because the problem is equivalent to finding collisions for a function $\mathbf{H}_1(x) | \mathbf{H}_2(x)$ that is compressing, which necessarily exist. \square

The birthday attack gives a $2^{\min(n+s_1, n+s_2)/2}$ upper bound on the security of the combined hash function. Balancing the digest lengths with $s_1 = s_2 = n/2$, leads to a maximum collision security of at most $2^{n/4}$. Proving a lower bound, on the other hand, remains an interesting open problem. We formulate a conjecture towards proving this next.

Conjecture 1. The multi-instance 2-element set-intersection problem over Bernoulli sets in a universe of size N with $p = q = 1/\sqrt{N}$ and \sqrt{N} sets for each party has communication complexity

$$D_{\mu(p,q)}^\varepsilon((\sqrt{N}, \sqrt{N})\text{-2INT}) \geq \tilde{\Omega}(N^{1/4})$$

for a sufficiently small protocol error ε and where $\tilde{\Omega}$ hides logarithmic factors.

We note that a lower bound for protocols with a sufficiently small error would be sufficient for feasibility results as collision resistance can also be amplified in a black-box way [8].

6 The Cascade Combiner

We now look at the security of the cascade combiner

$$\mathbf{C}^{\mathbf{H}_1, \mathbf{H}_2}(x) := \mathbf{H}_2(\mathbf{H}_1(x))$$

in the 2-BRO model, where $\mathbf{H}_1 \in \text{Fun}[n, n + s_1]$ and $\mathbf{H}_2 \in \text{Fun}[n + s_1, n + s_1 + s_2]$. We will prove one-way security and pseudorandomness based on set-intersection and set-disjointness respectively, and collision resistance based on a variant finding two intersecting points given multiple instances for one party and a single set for the other.

6.1 One-way security

Similarly to the concatenation combiner, we can reduce the random preimage resistance (rPre) security of the cascade combiner to set-intersection.

Theorem 7. *Let $H_1 \in \text{Fun}[n, n + s_1]$ and $H_2 \in \text{Fun}[n + s_1, n + s_1 + s_2]$ and $C^{H_1, H_2}(x) := H_2(H_1(x))$. Then for large enough n and any adversary \mathcal{A} against the rPre security of C^{H_1, H_2} in the 2-BRO model there is a 2-party protocol π against set-intersection with $\mu := \mu(p, q)$ where $p := 1/2^{s_1}$ and $q := 1/2^{n+s_1+s_2}$ and such that*

$$\text{Adv}_{C^{H_1, H_2}}^{\text{rpre}}(\mathcal{A}) \leq \text{Adv}_{\mu}^{\text{int}}(\pi) + \sqrt{n}2^{-n/2}(1+2^{s_2-s_1}) \quad \text{and} \quad D_{\mu}(\pi) \leq Q(\mathcal{A}) + 3n + s_1 + s_2.$$

Proof. We follow a strategy similar to the reductions in Section 5. Given a random $y^* \in \{0, 1\}^{n+s_1+s_2}$ the task of the adversary \mathcal{A} against rPre security of C^{H_1, H_2} is to find a z such that $C^{H_1, H_2}(z) = y^*$. With such a z , one can then also compute $x := H_1(z)$ and conclude that $x \in I \cap T$ where

$$I := H_1(\{0, 1\}^n) \quad \text{and} \quad T := H_2^{-1}(y^*)$$

with $I, T \subseteq \{0, 1\}^{n+s_1}$. The set T is Bernoulli with parameter $\Pr[y \in T] = 1/2^{n+s_1+s_2}$. Although set I appears to be Bernoulli,

$$\Pr[x \in I] = 1 - \Pr[\forall z : H_1(z) \neq x] = 1 - (1 - 1/2^{n+s_1})^{2^n}$$

it is not, since these probabilities are not independent for different values of x .

Our strategy to deal with this and ultimately construct a protocol π for solving set-intersection is to start with a Bernoulli set S (Alice's input), and program H_1 on all $x \in \{0, 1\}^n$ to values y that will be taken from S , but are also set to *collide* with the right probability. This will ensure that the image of H_1 contains most of S and is also distributed as the image of a random oracle.

We proceed as follows. Initially the set of assigned domain points X and assigned co-domain points Y are empty. We then iterate through $x \in \{0, 1\}^n$ in a random order. A bit b decides at each iteration decides if the hash value y for x should collide with a previously assigned value or not. If so, we sample y from the set of already assigned values Y . Otherwise, y should be a non-colliding value and we sample it from S if S is non-empty (and remove y from S), or otherwise we sample it outside the already assigned points Y . The pseudo-code for this algorithm, which we call **HashSam**, is shown in Figure 5.

Setting $m := n + s_1$, we now need to check that (1) the returned hash function H_1 is distributed as a random oracle $\{0, 1\}^n \rightarrow \{0, 1\}^m$ when S is Bernoulli with parameter $p = 1/2^{s_1}$, and (2) if $x \in H_1(\{0, 1\}^n) \cap H_2^{-1}(y^*)$, then we also have that $x \in S \cap T$ with good probability.

We first prove (1). The intuition is that the algorithm treats all inputs and outputs in a uniform way, and hence no particular values are special. Formally, let x^* and y^* be any fixed values. We show that $\Pr[H_1(x^*) = y^*] = 1/2^m$, even given the previously assigned values. We use a subscript i to denote the values

<p>Algo. HashSam(S)</p> <hr/> <p>$X \leftarrow \emptyset; Y \leftarrow \emptyset$</p> <p>for $i = 1, \dots, 2^n$ do</p> <p style="padding-left: 20px;">$x \leftarrow \{0, 1\}^n \setminus X; X \leftarrow X \cup \{x\}$</p> <p style="padding-left: 20px;">$b \leftarrow \text{Ber}(Y /2^m)$</p> <p style="padding-left: 20px;">if $b = 1$ then $y \leftarrow Y$</p> <p style="padding-left: 20px;">if $b = 0 \wedge S = \emptyset$ then</p> <p style="padding-left: 40px;">$y \leftarrow \{0, 1\}^m \setminus Y; Y \leftarrow Y \cup \{y\}$</p> <p style="padding-left: 20px;">if $b = 0 \wedge S \neq \emptyset$ then</p> <p style="padding-left: 40px;">$y \leftarrow S; Y \leftarrow Y \cup \{y\}; S \leftarrow S \setminus \{y\}$</p> <p style="padding-left: 20px;">$H_1 \leftarrow H_1 : [x \mapsto y]$</p> <p>return H_1</p>

Fig. 5: Hash sampler centered around a Bernoulli set S .

of various variables in the i -th iteration. Looking at different execution branches of the algorithm we can calculate $\Pr[y_i = y^* | x_i = x^*, Y_i, X_i]$ as

$$\Pr[b_i = 1] \Pr[y^* \in Y_i] \frac{1}{|Y_i|} + \Pr[b_i = 0] \left(\Pr[S_i = \emptyset] \Pr[y^* \notin Y_i] \frac{1}{2^m - |Y_i|} + \Pr[S_i \neq \emptyset] \Pr[y^* \in S_i] \frac{1}{|S_i|} \right).$$

Letting $\theta_i := \Pr[S_i = \emptyset]$ we can simplify to

$$\frac{|Y_i|}{2^m} \frac{|Y_i|}{2^m} \frac{1}{|Y_i|} + \left(1 - \frac{|Y_i|}{2^m}\right) \left(\theta_i \left(1 - \frac{|Y_i|}{2^m}\right) \frac{1}{2^m - |Y_i|} + (1 - \theta_i) \frac{|S_i|}{2^m} \frac{1}{|S_i|} \right) = \frac{1}{2^m}.$$

Note we have used the fact that S_i is a Bernoulli set in $\Pr[y^* \in S_i] = \frac{|S_i|}{2^m}$. Hence

$$\Pr[H_1(x^*) = y^* | Y_i, X_i] = \sum_{i=1}^{2^n} \Pr[y_i = y^* | x_i = x^*, Y_i, X_i] \Pr[x_i = x^*] = \frac{1}{2^m}.$$

Therefore the probability of sampling any given hash function is $(1/2^m)^{2^n}$, as required.

Let us now consider (2). When $I \subseteq S$, any solution with respect to I is also one with respect to S (that is, solutions are not lost). Hence we only look at the case $S \subseteq I$ and bound $|I \setminus S| = |I| - |S|$. Since $|I| \leq 2^n$ and $\mathbb{E}[|S|] = 2^{n+s_1}/2^{s_1} = 2^n$, we get that for any t

$$\Pr[|I| - |S| > t] \leq \Pr[2^n - |S| > t] = \Pr[\mathbb{E}[|S|] - |S| > t].$$

Applying the Chernoff bounds we obtain

$$\Pr \left[\mathbb{E}[|S|] - |S| > t \mathbb{E}[|S|] \right] \leq e^{-\frac{t^2}{2+t} \mathbb{E}[|S|]}.$$

Setting $t := \sqrt{n/2^n}$, we get with overwhelming probability that $|I \setminus S| \leq \sqrt{n}2^{n/2}$. Hence $T \cap (I \setminus S)$ will be non-empty with negligible probability $\sqrt{n}2^{-n/2-s_1-s_2}$, in which case if $x \in I \cap T \implies x \in S \cap T$. \square

If $H_1 \in \text{Fun}[n, (2 + \epsilon)n]$ and $H_2 \in \text{Fun}[(2 + \epsilon)n, (1 + \epsilon)n]$, we have a reduction to set-intersection with parameters $N = 2^{(2+\epsilon)n}$, $p = 1/2^{(1+\epsilon)n}$, and $q = 1/2^{(1+\epsilon)n}$. Thus with notation as in the description of the feasible in Figure 3 we have that $\alpha = \beta = (1 + \epsilon)/(2 + \epsilon)$. As in Corollary 1 we would need The point (α, β) lies in the feasible region for $1/3 < (1 + \epsilon)/(2 + \epsilon) < 1/2$, which means $-1/2 < \epsilon < 0$. Since the combined function is compressing (with $\gamma = 1/(1 + \epsilon) > 1$) and $p_o \approx 1 - e^{-2^{-\epsilon n}}$ is negligible, the image uniformity bound is negligible and hence, similarly to Corollary 1 we get strong OW security.

6.2 PRG and CR security

We briefly outline how to treat the PRG security and collision resistance of cascade. We omit the proofs as the techniques and proof structures are similar to our other results above.

PRG SECURITY. We can prove an analogous result for the oblivious PRG security of the cascade construction. Its reduction is identical to that for rPre security given above, except that the underlying assumption is set-disjointness. Setting $s_1 = 2n$ (H_1 is length doubling) and $s_2 = -2n + 1$ (H_2 compresses by almost a factor of 3) leads to a reduction to an instance of set-intersection with parameters $N = 2^{3n}$, $p = 1/2^{2n}$, and $q = 1/2^{n+1}$. In this case $\delta = 2$ and $p_o = 1 - e^{-1/2}$. With these parameters we can carry out an analysis similar to Corollary 2: We set the error $\varepsilon = 0.065$ which is smaller than $(\delta - 1)p_o/(4 + \delta) < 0.0656$ as required in Theorem 2 for an exponential number of queries. The combined hash function maps n bits to $n + 1$ bits and hence $C = 1/2$. Furthermore, p_o is negligible as a function from n bits to $3n$ bits is injective with overwhelming probability. Thus we can apply Lemma 3 with $s = 1$ as in Corollary 2 to get a weak PRG.

COLLISION RESISTANCE. We can treat the collision resistance of cascade similarly. The difference is that in the reduction Alice will use the HashSam algorithm in Figure 5 to adapt a (single) Bernoulli set S that she holds to a hash image set I . On the other hand, Bob uses the ReDist algorithm in Figure 4 to redistribute elements in multiple Bernoulli sets that he holds so that they form a partition of the entire domain of H_2 . The rest of the proof, which is included in the full version [4, Appendix E], proceeds similarly to Lemma 6. For setting parameters, observe that any collision for H_1 is necessarily a collision for $H_2(H_1(\cdot))$. Since collisions for H_1 can be easily found using BD_1 , we need H_1 to be injective. For example, $s_1 = 2n$ (co-domain points are $3n$ bits) would lead to an injective H_1 with overwhelming probability.

Acknowledgments. We thank Marc Fischlin for participating in the early stages of this work. We also thank the CRYPTO'18 (sub)reviewers for their

valuable comments. Bauer was supported by the French ANR Project ANR-16-CE39-0002 EfTrEC. Farshim was supported by the European Research Council under the European Community’s Seventh Framework Programme (FP7/2007-2013 Grant Agreement no. 339563 - CryptoCloud). Mazaheri was supported by the German Federal Ministry of Education and Research (BMBF) and by the Hessian State Ministry for Higher Education, Research and the Arts, within CRISP.

References

1. L. Babai, P. Frankl, and J. Simon. Complexity classes in communication complexity theory (preliminary version). In *27th FOCS*, pages 337–347, 1986.
2. Z. Bar-Yossef, T. S. Jayram, R. Kumar, and D. Sivakumar. An information statistics approach to data stream and communication complexity. In *43rd FOCS*, pages 209–218, 2002.
3. B. Barak, M. Braverman, X. Chen, and A. Rao. How to compress interactive communication. In *42nd ACM STOC*, pages 67–76, 2010.
4. B. Bauer, P. Farshim, and S. Mazaheri. Combiners for backdoored random oracles. Cryptology ePrint Archive, 2018.
5. M. Bellare and P. Rogaway. Random oracles are practical: A paradigm for designing efficient protocols. In *ACM CCS 93*, pages 62–73, 1993.
6. D. J. Bernstein, T. Lange, and R. Niederhagen. Dual EC: A standardized back door. Cryptology ePrint Archive, Report 2015/767, 2015. <http://eprint.iacr.org/2015/767>.
7. D. Boneh and X. Boyen. On the impossibility of efficiently combining collision resistant hash functions. In *CRYPTO 2006*, pages 570–583, 2006.
8. R. Canetti, R. L. Rivest, M. Sudan, L. Trevisan, S. P. Vadhan, and H. Wee. Amplifying collision resistance: A complexity-theoretic treatment. In *CRYPTO 2007*, pages 264–283, 2007.
9. A. Chattopadhyay and T. Pitassi. The story of set disjointness. *SIGACT News*, 41(3):59–85, 2010.
10. S. Checkoway, J. Maskiewicz, C. Garman, J. Fried, S. Cohney, M. Green, N. Heninger, R.-P. Weinmann, E. Rescorla, and H. Shacham. A systematic analysis of the juniper dual EC incident. In *ACM CCS 16*, pages 468–479, 2016.
11. S. Checkoway, R. Niederhagen, A. Everspaugh, M. Green, T. Lange, T. Ristenpart, D. J. Bernstein, J. Maskiewicz, H. Shacham, and M. Fredrikson. On the practical exploitability of dual EC in TLS implementations. In *23rd USENIX Security Symposium (USENIX Security 14)*, pages 319–335, 2014.
12. S. Coretti, Y. Dodis, S. Guo, and J. Steinberger. Random oracles and non-uniformity. Cryptology ePrint Archive, Report 2017/937, 2017. <http://eprint.iacr.org/2017/937>.
13. I. Dinur. New attacks on the concatenation and XOR hash combiners. In *EUROCRYPT 2016, Part I*, pages 484–508, 2016.
14. Y. Dodis, S. Guo, and J. Katz. Fixing cracks in the concrete: Random oracles with auxiliary input, revisited. In *EUROCRYPT 2017, Part II*, pages 473–495, 2017.
15. A. Fiat and A. Shamir. How to prove yourself: Practical solutions to identification and signature problems. In *CRYPTO’86*, pages 186–194, 1987.
16. M. Fischlin and A. Lehmann. Security-amplifying combiners for collision-resistant hash functions. In *CRYPTO 2007*, pages 224–243, 2007.

17. M. Fischlin, A. Lehmann, and K. Pietrzak. Robust multi-property combiners for hash functions. *Journal of Cryptology*, 27(3):397–428, 2014.
18. O. Goldreich. *Foundations of Cryptography: Basic Tools*, volume 1. Cambridge University Press, Cambridge, UK, 2001.
19. V. Guruswami and M. Cheraghchi. Set disjointness lower bound via product distribution. Scribes for Information theory and its applications in theory of computation, 2013. <http://www.cs.cmu.edu/~venkatg/teaching/ITCS-spr2013/>.
20. J. J. Hoch and A. Shamir. On the strength of the concatenated hash combiner when all the hash functions are weak. In *ICALP 2008, Part II*, pages 616–630, 2008.
21. A. Joux. Multicollisions in iterated hash functions. Application to cascaded constructions. In *CRYPTO 2004*, pages 306–316, 2004.
22. J. Katz, S. Lucks, and A. Thiruvengadam. Hash functions from defective ideal ciphers. In *CT-RSA 2015*, pages 273–290, 2015.
23. A. Kawachi, A. Numayama, K. Tanaka, and K. Xagawa. Security of encryption schemes in weakened random oracle models. In *PKC 2010*, pages 403–419, 2010.
24. E. Kushilevitz and N. Nisan. *Communication complexity*. Cambridge University Press, 1997.
25. A. Lehmann. *On the Security of Hash Function Combiners*. PhD thesis, TU Darmstadt, 2010.
26. G. Leurent and L. Wang. The sum can be weaker than each part. In *EUROCRYPT 2015, Part I*, pages 345–367, 2015.
27. M. Liskov. Constructing an ideal hash function from weak ideal compression functions. In *SAC 2006*, pages 358–375, 2007.
28. S. Lucks. A failure-friendly design principle for hash functions. In *ASIACRYPT 2005*, pages 474–494, 2005.
29. U. M. Maurer and S. Tessaro. A hardcore lemma for computational indistinguishability: Security amplification for arbitrarily weak PRGs with optimal stretch. In *TCC 2010*, pages 237–254, 2010.
30. F. Mendel, C. Rechberger, and M. Schl affer. MD5 is weaker than weak: Attacks on concatenated combiners. In *ASIACRYPT 2009*, pages 144–161, 2009.
31. A. Mittelbach. Cryptophia’s short combiner for collision-resistant hash functions. In *ACNS 13*, pages 136–153, 2013.
32. D. Moshkovitz and B. Barak. Communication complexity. Scribes for Advanced Complexity Theory, 2012. <https://people.csail.mit.edu/dmoshkov/courses/adv-comp/>.
33. A. Numayama, T. Ishiki, and K. Tanaka. Security of digital signature schemes in weakened random oracle models. In *PKC 2008*, pages 268–287, 2008.
34. O. Reingold, L. Trevisan, and S. P. Vadhan. Notions of reducibility between cryptographic primitives. In *TCC 2004*, pages 1–20, 2004.
35. P. Rogaway and T. Shrimpton. Cryptographic hash-function basics: Definitions, implications, and separations for preimage resistance, second-preimage resistance, and collision resistance. In *FSE 2004*, pages 371–388, 2004.
36. D. R. Simon. Finding collisions on a one-way street: Can secure hash functions be based on general assumptions? In *EUROCRYPT’98*, pages 334–345, 1998.
37. D. Unruh. Random oracles and auxiliary input. In *CRYPTO 2007*, pages 205–223, 2007.
38. A. C.-C. Yao. Some complexity questions related to distributive computing (preliminary report). In *Proceedings of the Eleventh Annual ACM Symposium on Theory of Computing*, pages 209–213, 1979.