



# On the Compressibility of $\mathcal{NP}$ Instances and Cryptographic Applications<sup>‡</sup>

Danny Harnik\*

Moni Naor<sup>†</sup>

## Abstract

We initiate the study of the compressibility of  $\mathcal{NP}$  problems. We consider  $\mathcal{NP}$  problems that have long instances but relatively short witnesses. The question is, can one efficiently compress an instance and store a shorter representation that maintains the information of whether the original input is in the language or not. We want the length of the compressed instance to be polynomial in the length of the *witness* rather than the length of original input. Such compression enables to succinctly store instances until a future setting will allow solving them, either via a technological or algorithmic breakthrough or simply until enough time has elapsed.

We give a new classification of  $\mathcal{NP}$  with respect to compression. This classification forms a stratification of  $\mathcal{NP}$  that we call the  $\mathcal{VC}$  hierarchy. The hierarchy is based on a new type of reduction called  $W$ -reduction and there are compression-complete problems for each class.

Our motivation for studying this issue stem from the vast cryptographic implications compressibility has. For example, suppose that SAT is compressible, that is there exist a polynomial  $p(\cdot, \cdot)$  so that given a formula consisting of  $m$  clauses over  $n$  variables it is possible to come up with an equivalent (w.r.t satisfiability) formula of size at most  $p(n, \log m)$ . Then if the reduction is what we call *witness retrievable* we provide a construction of an Oblivious Transfer Protocol from *any* one-way function. Using the terminology of Impagliazzo [Imp95], this implies that  $\text{Minicrypt} = \text{Cryptomania}$ . Other implications of SAT compressibility (without the witness retrievability) are: (i) the construction of collision resistant hash function from *any* one-way function and (ii) a cryptanalytic result concerning the limitation of everlasting security in the bounded storage model when mixed with (time) complexity based cryptography.

---

\*Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot 76100 Israel. Email: danny.harnik@weizmann.ac.il.

<sup>†</sup>Incumbent of the Judith Kleeman Professorial Chair, Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot 76100 Israel. Email: moni.naor@weizmann.ac.il.

<sup>‡</sup>Research supported by a grant from the Israel Science Foundation.

# 1 Introduction

In order to deal with difficult computational problems several well established options were developed, including: approximation algorithms, subexponential algorithms, parametric complexity and limited non-determinism. In this paper we explore our favorite approach for dealing with problems: postpone them (hopefully, without cluttering our desk). We initiate the study of the compressibility of NP problems for their resolution in some future setting. Rather than solving a given instance, we ask whether a shorter instance with the same solution can be found efficiently. Specifically, we consider  $\mathcal{NP}$  problems that have long instances but relatively short witnesses. An  $\mathcal{NP}$  language  $L$  is defined by an efficiently computable relation  $R_L$  such that an input (or instance)  $x$  is in  $L$  if and only if there exists a witness  $w$  such that  $R_L(x, w) = 1$ . Throughout the paper, an  $\mathcal{NP}$  instance is characterized by two parameters  $m$  and  $n$ : The length of the instance  $x$  is denoted by  $m$  and the length of the witness  $w$  is denoted by  $n$ . The problems of interest are those having short witnesses, i.e.  $n \ll m$ . Traditionally, the study of  $\mathcal{NP}$  languages evolves around the ability or inability to efficiently decide if an instance is in the language or not, or to find a witness  $w$  for an instance  $x$  within polynomial time. We introduce the question of compressibility of such instances.

**Compressing SAT Instances:** To illustrate the relevant setting, we use the well known example of SAT. An instance  $\Phi$  for SAT consists of a CNF formula over  $n$  variables and we define that  $\Phi \in SAT$  if there exists an assignment to the  $n$  variables that satisfies all the clauses of  $\Phi$ . The question of compressibility of SAT is the following:

## Example 1.1 (Compression of SAT instances)

*Does there exist an efficient algorithm and a polynomial  $p(\cdot, \cdot)$  with the following input and output:*

**Input:** A CNF formula  $\Phi$  with  $m$  clauses over  $n$  variables (we are interested in  $m \gg n$ ).

**Output:** A formula  $\Psi$  of size  $p(n, \log m)$  such that  $\Psi$  is satisfiable if and only if  $\Phi$  is satisfiable.

The idea is that the length of  $\Psi$  should be essentially unrelated to the original length  $m$ , but rather to the number of variables (or in other words, to the size of the witness). Typically, we think of the parameters  $m$  and  $n$  as related by some function, and it is instructive (but not essential) to think of  $m$  as larger than any polynomial in  $n$ . So potentially, the length of  $\Psi$  can be significantly shorter than that of  $\Phi$ .<sup>1</sup>

In general, one cannot expect to compress all the formulas, or else we would have an efficient algorithm for all  $\mathcal{NP}$  problems.<sup>2</sup> However, once we introduce the setting of a shorter witness, then compression becomes plausible. Note that if  $\mathcal{P} = \mathcal{NP}$  and we actually know the algorithm for SAT then clearly compression is trivial, simply by solving the satisfiability of  $\Phi$  and outputting 1 if  $\Phi \in SAT$  and 0 otherwise.

**Motivation for Compression:** Compressing for the future is an appealing notion for various settings. There are numerous plausible scenarios that will give us more power to solve problems in the future. We could potentially find out that  $\mathcal{P} = \mathcal{NP}$  and solve all our  $\mathcal{NP}$  problems then. We may have faster computers or better means of computing such as quantum computers or any other physical method for solving problems (see Aaronson [Aar05] for a list of suggestions). Above all, the future entails lots and lots of time, a resource that the present is usually short of. Saving the problems of today as they are presented is wasteful, and compression of problems will allow us to store a far greater number of problems for better days.

---

<sup>1</sup>Note, that since our requirement for compression is only relevant for problems where  $m \gg n$ , then an  $\mathcal{NP}$ -complete problem such as 3-SAT (where all clauses have exactly 3 literals) is irrelevant for compression as in such formulas  $m$  is already at most  $O(n^3)$ .

<sup>2</sup>Suppose that every formula can be compressed by a single bit, then sequentially reapplying compression to the input will result in a very short formula that may be solved by brute enumeration.

Our interest in studying the issue of compression stems from the vast cryptographic implications of compressibility. We demonstrate three questions in cryptography that compression algorithms would resolve (see Section 1.3). We are confident that the compression of problems implies further applications both within and outside of cryptography. For example, Dubrov and Ishai [DI06], in subsequent work, show the relevance of the notion of compression to derandomization (see related work, Section 1.4). The concept of compression of problems is also interesting beyond the confines of  $\mathcal{NP}$  problems, and makes sense in any setting where the compression requires much less resources than the actual solution of the problem.

## 1.1 Compression of NP instances

We define the notion of compression with respect to an  $\mathcal{NP}$  language. For simplicity, we assume that an input to an  $\mathcal{NP}$  language  $L$  includes an encoding of the parameter  $n$  that upper bounds the length of a potential witness.<sup>3</sup> We also associate with  $L$  a specific  $\mathcal{NP}$  relation  $R_L$  that defines it (as mentioned above). We note that once the parameters  $m$  and  $n$  are explicit, it is in most cases immaterial what specific relation defines the language and the properties we discuss (such as compressibility) are properties of the language at hand (unless stated otherwise).

In essence, a compression algorithm is a specialized Karp-reduction that also reduces the length of the instance. In the following definition we also introduce an additional property that is of interest for some of the applications. This is called *witness retrieval* and essentially means that knowledge of a witness to the input instance implies knowledge of a witness to the compressed instance.<sup>4</sup>

**Definition 1.2 (Compression Algorithm for  $\mathcal{NP}$  Instances)** *Let  $L$  be an  $\mathcal{NP}$  language. A compression algorithm for  $L$  is a probabilistic polynomial time machine  $Z$  along with a language  $L'$  in  $\mathcal{NP}$  (or more accurately in  $\mathcal{NP}(poly(m))$ )<sup>5</sup> and a polynomial  $p(\cdot, \cdot)$  such that for all large enough  $m$ :*

1. For all  $x \in \{0, 1\}^m$  with parameter  $n$  the length of  $Z(x)$  is at most  $p(n, \log m)$ .
2.  $Z(x) \in L'$  if and only if  $x \in L$

*A compression algorithm is witness retrievable with respect to  $R_L$  if there exists a probabilistic polynomial time machine  $W$  such that if  $w_x$  is a witness for  $x \in L$  then  $w_y = W(w_x, Z(x))$  is a witness for  $Z(x) \in L'$ .*

*We allow a negligible error in the success of  $Z$  and  $W$  (where probability is over the internal randomness of  $Z$  and  $W$ ).*

The paper consists of two parts: *Part I* is a study of the concept of compression of  $\mathcal{NP}$  instances from a complexity point of view. *Part II* introduces the cryptographic applications of compression algorithms.

**How much to compress:** Definition 1.2 (of compression algorithms) requires a very strong compression, asking that the length of the compression is polynomial in  $n$  and  $\log m$ . For the purposes of part I of the paper (the complexity study), it is essential that the compression is at least sub-polynomial in  $m$  in order to ensure that the reductions defined with respect to compressibility (See Section 2.2) do compose.<sup>6</sup> Furthermore,

<sup>3</sup>Typically, the parameter  $n$  is indeed part of the description of the problem (e.g. for Clique, SAT, Sparse Subset Sum and others).

<sup>4</sup>The reverse direction, where one wants to extract a witness to the input instance from the compressed instance is considered in the Section 2.7 on the compression of search problems.

<sup>5</sup>By  $\mathcal{NP}(poly(m))$  we mean in nondeterministic-time  $poly(m)$  (that is, verifiable in time  $poly(m)$  when given a non-deterministic hint).

<sup>6</sup>For clarity we choose a polynomial in  $\log m$ , although this may be replaced by any sub-polynomial function  $m'(\cdot)$  (a function such that for large enough  $m$  for any polynomial  $q(\cdot)$  we have  $m'(m) < q(m)$ ).

for part II (the applications) this definition may be strongly relaxed, where even a compression to  $m^{1-\varepsilon}$  for some constant  $\varepsilon$  suffices for some applications.

**The Complexity of  $L'$ :** Another requirement of Definition 1.2 is that the language  $L'$  be in  $\mathcal{NP}(poly(m))$ . In general, this requirement may also be relaxed and the result still be meaningful for some applications. In particular, we do not need to put a bound on the complexity of  $L'$ , but only require that there is enough information in  $Z(x)$  to determine whether  $x \in L$  or not. One case where we use a definition with unbounded extraction is the compression of search problems in Section 2.7. It should be noted however that the requirement for  $L'$  to be in  $\mathcal{NP}(poly(m))$  is necessary in order for the witness retrievability property to be meaningful. Moreover, in some cases it is natural to further restrict  $L'$  to actually be in  $\mathcal{NP}$  (that is in  $\mathcal{NP}(poly(n, \log m))$ ). For instance, this is the case in the definition of compression of SAT (Example 1.1). Finally, note that if the compression has *zero error probability*, then  $L'$  must be in  $\mathcal{NP}(poly(m))$  simply by the definition of compression.<sup>7</sup>

## 1.2 Part I: Classifying $\mathcal{NP}$ Problems with Respect to Compression.

We are interested in figuring out which  $\mathcal{NP}$  languages are compressible and, in particular, whether important languages such as SAT and Clique are compressible. For starters, we demonstrate some non-trivial languages that do admit compression: we show compression for the well known  $\mathcal{NP}$ -complete problem of vertex-cover and for another  $\mathcal{NP}$ -complete language known as minimum-fill-in. We also show compression of a language consisting of strings that are the output of a cryptographic pseudorandom generator as well as compression for the promise problem GapSAT.<sup>8</sup> However, these examples are limited and do not seem to give us compression for all  $\mathcal{NP}$ . Moreover, it becomes clear that the traditional notions of reductions and completeness in  $\mathcal{NP}$  do not apply for the case of compression (i.e., the compression of an  $\mathcal{NP}$ -complete language does not immediately imply compression for all of  $\mathcal{NP}$ ). Neither are different notions of reduction with respect to other aspects of computation, such as reductions with respect to approximation algorithms or subexponential algorithms (see for example [Pap94]) and for parameterized complexity (see [DF99] and further discussion in Section 1.4 on related work).

We introduce  $W$ -reductions that arise from the study of compression. These are reductions that address the length of the witness in addition to membership in an  $\mathcal{NP}$  language. Following the definition of  $W$ -reductions we define also the matching notion of compression-complete languages for a class. We then introduce a classification of  $\mathcal{NP}$  problems with respect to compression. The classification presents a structured hierarchy of  $\mathcal{NP}$  problems, that is surprisingly different from the traditional view and closer in nature to the  $W$  hierarchy of parameterized complexity (see [DF99]). We call our hierarchy  $\mathcal{VC}$ , short for “verification classes”, since the classification is closely related to the verification algorithm of  $\mathcal{NP}$  languages when allowed a preprocessing stage. We discuss some of the more interesting classes in the  $\mathcal{VC}$  hierarchy, mention some compression-complete problems and classify some central  $\mathcal{NP}$  problems.

In addition, we study the compression of  $\mathcal{NP}$  search problems. That is, compressing an instance in a way that maintains all the information about a witness for the problem. We show that the compression of a class of decision problems also implies compression for the corresponding search problems.

<sup>7</sup>Suppose that there exists a compression algorithm  $Z$  for  $L$  then define  $L'$  to be the language of all  $Z(x)$  such that  $x \in L$ . Then, for every  $y \in L'$  a verification algorithm takes as a nondeterministic witness a value  $x$ , a witness to  $x \in L$  along with randomness for the compression algorithm and verifies that indeed  $y = Z(x)$ . Thus if  $Z$  never introduces an error then  $L'$  is in  $\mathcal{NP}(poly(m))$ .

<sup>8</sup>I.e. a promise problem were either the formula is satisfiable or every assignment does not satisfy a relatively large number of clauses.

### 1.3 Part II: Implications to Cryptography

We provide some implications of compressibility to cryptography. The implications described are of contrasting flavors. On the one hand we show constructions of cryptographic primitives using compression algorithms, while on the other hand we show a cryptanalysis using compression algorithms (or alternatively, this can be considered as an application of incompressibility of languages). For simplicity we provide the implication with respect to the compression of SAT. We note however, that the same statements can actually be made with compression of languages from the class  $\mathcal{VC}_{OR}$  (see Definition 2.17). This class is the lowest class in our  $\mathcal{VC}$  hierarchy, and potentially easier to compress than SAT.

**On the existence of Minicrypt:** Impagliazzo [Imp95] summarizes five possibilities for how the world may look like based on different computational assumptions. The two top worlds are `Minicrypt`, where one-way functions exist but oblivious transfer protocols do not exist (in this world some interesting cryptographic applications are possible, and in particular *shared* key cryptography) and `Cryptomania` where Oblivious Transfer protocols do exist (and hence also a wide range of cryptographic applications like secure computation and *public* key cryptography). We show that if there exists a compression algorithm for SAT that is also witness retrievable, then `Minicrypt`=`Cryptomania`:

**Theorem 1.3** *If there exists a witness retrievable compression algorithm for SAT, then there exists an Oblivious Transfer protocol based on any one-way function.*

This theorem is proved by first constructing a private information retrieval (PIR) protocol from any one-way function and then applying the reduction of Di Crescenzo, Malkin and Ostrovsky [DMO00] from OT to PIR. As immediate corollaries we also get constructions for general secure computation protocols (e.g. [Yao82, GMW87, GV87, Kil88]) and secret key agreement from one-way functions. All of the above are *not* known to be equivalent to the existence of one-way functions. Moreover, Impagliazzo and Rudich [IR89] prove that key agreement protocols (and hence also OT) cannot be constructed from any one-way function using black-box reductions. This does not imply that witness retrievable compression of SAT is impossible, since the construction of OT in Theorem 3.1 is inherently non-black-box and uses the program of the one-way function via Cook's Theorem [Coo71]. Whether OT can be constructed from any one-way function is a major open problem in cryptography and as such may be viewed as an indication to the difficulty of resolving the question of finding witness retrievable compression algorithms for SAT.

**On Collision Resistant Hash from any One-Way Function:** Practically the same construction of PIR from Theorem 3.1 but *without the witness retrievability* can be used, via a result of Ishai, Kushilevitz and Ostrovsky [IKO05], to construct collision resistant hash functions (CRH). Thus we get the following:

**Theorem 1.4** *If there exists an errorless<sup>9</sup> compression algorithm for SAT then there exists a construction of collision resistant hash functions based on any one-way function.*

This task as well was shown to be impossible by black-box reductions (see Simon [Sim98]). Another interesting corollary of this result is a construction of statistically hiding bit commitment from any one-way function, which is currently an open problem ([NOVY98, HHK<sup>+</sup>05] show such constructions based on one-way functions with a specific structure).

---

<sup>9</sup>The construction of CRH requires that the error probability of compression algorithm will be zero. This can be slightly relaxed to an error that is exponentially small in  $m$  (rather than  $n$ ).

**On Everlasting Security and the Hybrid Bounded Storage Model:** The *bounded storage model* (BSM) of Maurer [Mau92] provides the setting for the appealing notion of *everlasting security* [ADR02, DR02]. Loosely speaking, two parties, Alice and Bob, that share a secret key in advance, may use the BSM to encrypt messages in a way that the messages remain secure against a computationally unbounded adversary, even if the shared secret key is eventually revealed.

However, if the parties do not meet in advance to agree on a secret key then everlasting security requires high storage requirements from Alice and Bob [DM04a], rendering encryption in this model less appealing. Hoping to overcome this, it was suggested to combine the BSM with computational assumptions (what is called here the hybrid BSM). In particular, to run a computational key agreement protocol in order to agree on a shared secret key, and then run one of the existing BSM schemes. Dziembowski and Maurer [DM04a] showed that this idea does not necessarily work in all cases, by showing an attack on a protocol consisting of the combination of a specific (artificial) computational key agreement protocol with a specific BSM encryption scheme.

We show that compression of  $\mathcal{NP}$  instances can be used to attack *all* hybrid BSM schemes. Or in other words, if a compression of SAT exists, then the hybrid BSM is no more powerful than the standard BSM.

One interpretation of this result is that in order to prove everlasting security for a hybrid BSM scheme, without further conditions, one must prove that there exists no compression algorithm for SAT. Alternatively, as a relaxation, one should come up with a reasonable incompressibility assumption regarding the resulting formulae. Note however that a straightforward assumption of the form “this distribution on SAT formulae is incompressible” is not efficiently falsifiable, in the sense of Naor [Nao03], that is, it is not clear how to set up a challenge that can be solved in case the assumption is false.

ON RANDOM ORACLES: We show (in [HN05]) that if all parties are given access to a random oracle, then there actually exists everlasting security in the hybrid BSM without an initial key and with low storage requirements from Alice and Bob<sup>10</sup>. Therefore, finding a compression algorithm for SAT would show an example of a task that is simple with random oracles but altogether impossible without them. This is stronger than previous results (such as [CGH04, GK03, MRH04]) that show a specific protocol that becomes insecure if the random oracle is replaced by a function with a small representation. This would constitute an argument against relying (blindly) on random oracles to determine whether a task is feasible at all.

## 1.4 Related Work

The relationship between compression and complexity in general is a topic that has been investigated since the early days of Complexity Theory (i.e. Kolmogorov Complexity [LV97]). However, the feature that we are introducing in this work is compressibility with respect to the *solution* (witness) rather than the instance. The goal of maintaining the solution differs our work from a line of seemingly related works about notions of compression ([DLN96, TVZ04, Wee04] to name a few), all of which aim at eventually retrieving the input of the compression algorithm.

There are several examples of other relaxations to solving  $\mathcal{NP}$  problems in polynomial time. Each of these relaxations yields a corresponding classification of  $\mathcal{NP}$ . These classifications, like ours, are subtle and usually turn out to be very different than that of the traditional  $\mathcal{NP}$  classification. For example, Papadimitriou and Yannakakis [PY88] introduce L-reductions and the classes MAX NP and MAX SNP, with respect to approximation algorithms. Impagliazzo, Paturi and Zane [IPZ98] define reductions with respect to solution in sub-exponential time.

Perhaps the most relevant classification to ours is that of parameterized complexity (see the monograph on this subject by Downey and Fellows [DF99]). Parameterized complexity studies the tractability of prob-

<sup>10</sup>This does not contradict the compressibility of SAT, since the cryptanalytic result is not black-box and assumes access to the full description of the programs of Alice and Bob. Thus this result is not preserved in the presence of a random oracle.

lems when one of the parameters is considered to be fixed or very small. This is relevant to compression since typically this parameter is related to the length of the witness. On the one hand, some (but not all) parameterized complexity algorithms yield natural compression algorithms (see examples and discussion in Section 2.1). In addition, some (but certainly not all) compression algorithms may imply a parameterized complexity algorithm. Also the  $W$ -hierarchy of parameterized complexity is reminiscent of the  $\mathcal{VC}$ -hierarchy (they are both defined by reduction to circuits of bounded depth). However, our study of compression yields quite a different classification. This is mainly because in parameterized complexity the witness length is taken to be very small and as such, there is no restriction on running in time that is exponential (or higher) in this parameter. In compression, on the other hand, the parameter (witness length) is usually of substantial size (even if much smaller than the instance length).

A related notion to parameterized complexity that is reminiscent of our work is *limited non-determinism*, which started with the work of Kintala and Fischer [KF80], see survey by Goldsmith, Levy and Mundheck [GLM96]. This was further studied by Papadimitriou and Yannakakis [PY96] who defined a few syntactic classes within the class of polylog non-determinism ( $LOGNP$  and  $LOGSNP$ ). The interesting point is that several natural problems are complete for these classes. The notion of reduction used is the usual polynomial reduction and hence the classifications arising from this study are very different from our  $\mathcal{VC}$  hierarchy.

In subsequent work, Dubrov and Ishai [DI06] discussed the compression of problems and showed that a certain incompressibility assumption has an application to derandomization. Specifically they construct a generator that fools procedures that use more randomness than their output length. Their work was mostly conducted independently of ours, following conversations regarding an early phase of our work. In addition, inspired by our CRH construction, they prove that any one-way permutation can either be used for the above mentioned derandomization, or else can be used to construct a weak version of CRH.<sup>11</sup>

**Paper organization:** Section 2 studies the compressibility of  $\mathcal{NP}$  problems and introduces the  $\mathcal{VC}$  classification with this respect. Part II of the paper (the cryptographic application) appears in Sections 3 and 4. Section 3 shows the applications to the construction of OT and CRH from any one-way functions, while Section 4 presents the implication to the hybrid BSM.

## 2 Part I: On the Compression of $\mathcal{NP}$ Instances

Attempting to compress  $\mathcal{NP}$  instances requires a different approach than solving  $\mathcal{NP}$  problems. Intuitively, a solution for compression might arise while trying to solve the problem. While a full solution of an  $\mathcal{NP}$  problem may take exponential time, it is plausible that the first polynomial number of steps leaves us without an explicit solution but with a smaller instance. Indeed, some algorithms in the parameterized complexity world work like this (see some examples in the next section). On the other hand, we allow the possibility that the compressed version is actually harder to solve (computational time-wise) than the original one (and may require a somewhat longer witness altogether).

### 2.1 Examples of Compression Algorithms for some Hard Problems

We start by showing three examples of compression algorithms for problems that are conjectured not to be in  $\mathcal{P}$ . Two of these examples are  $\mathcal{NP}$ -complete problems, while the third is taken from cryptography.

---

<sup>11</sup>This weak version of CRH (like the stronger common version) cannot be constructed from any one-way permutation by black-box reductions. (in [Sim98]).

**Vertex Cover:** The well studied  $\mathcal{NP}$ -complete problem of Vertex-Cover receives as input a graph  $G = (V, E)$  and asks whether there exists a subset of vertices  $S \subseteq V$  of size at most  $k$  such that for every edge  $(u, v) \in E$  either  $u$  or  $v$  are in  $S$ . The parameters are the instance length  $m$ , which is at most  $O(|E| \log |V|)$ , and the witness length  $n = k \log |V|$

**Claim 2.1** *There exists a witness retrievable compression algorithm for Vertex-Cover.*

**Proof:** We are following the parameterized complexity algorithm for vertex-cover (presented in [DF99] and attributed to Buss). If a vertex-cover  $S$  of size  $k$  exists, then any vertex of degree greater than  $k$  must be inside the set  $S$ . The compression algorithm simply identifies all such vertices and lists them in the cover, while removing all their outgoing edges (that do not need to be covered by other vertices). The graph left after this process has maximal degree  $k$ , and furthermore all edges have at least one end in the cover. Thus, if the original graph has a  $k$  vertex cover then the total number of edges left is at most  $k^2$  (at most  $k$  vertices in the cover with at most  $k$  edges each). If there are more than  $k^2$  edges then the answer to the problem is NO, otherwise, such a graph can be represented by the list of all edges, which takes  $k^2 \log k$  bits. The compression can be made witness retrievable since if we use the original labels of vertices to store the new graph, then the original cover is also a cover for the new compressed graph.  $\square$

It is interesting to note that we do not know of a compression algorithm for the Clique problem or the Dominating Set problem, which are strongly linked to the vertex-cover problem in the traditional study of  $\mathcal{NP}$ , and in fact, in Theorems 3.1, 3.3 and 4.4 we show strong implications of a compression algorithm for these languages.

**PRG-Output:** The following compression algorithm works for any sparse language (that is, a language that contains only a small fraction of all possible inputs). A specific and interesting example is defined next:

**Example 2.2 (PRG-Output)** *Let  $G$  be a pseudorandom generator stretching an  $n$  bit seed to an  $m$  bit output (with  $m \gg n$ ). Define the language PRG-output over inputs  $y \in \{0, 1\}^m$  as*

$$L_G = \{y \mid \text{there exists an } x \text{ s.t. } G(x) = y\}$$

The language PRG-output is hard to solve for random instances as long as the underlying PRG is secure. Yet it has a simple compression algorithm. Note that simply saving, say, the first  $2n$  bits of the instance  $y$  is insufficient because if  $y$  only differs from  $G(x)$  in one bit, then this bit may be anywhere in the  $m$  bits.

**Claim 2.3** *There exists a witness retrievable compression algorithm for PRG-output.*

**Proof:** Let  $\mathcal{H}$  be a family of almost pairwise independent hash functions from  $m$  bits to  $2n$  bits. The compression algorithm simply chooses a random  $h \in \mathcal{H}$  and outputs  $(h(z), h)$ . The new language is  $L'_G = \{(z, h) \mid \text{there exists an } x \text{ s.t. } h(G(x)) = z\}$ .

Naturally, if  $y \in L_G$  then also  $(h(y), h) \in L'_G$  with the same witness (and thus the witness retrievability). On the other hand, if  $y \notin L_G$  then by the properties of  $\mathcal{H}$ , for every seed  $x$  we have that  $\Pr_h[h(G(x)) = h(y)] < O(2^{-2n})$ , and by a union bound over all  $x \in \{0, 1\}^n$ , we get  $\Pr_h[h(y) \in L'_G] < O(2^{-n})$ . Finally, since there are almost pairwise independent hash functions whose description is of length  $O(n) + \log m$  (for example see [NN93]), then the algorithm is indeed compressing.  $\square$

**Minimum Fill-In:** The minimum fill-in problem is an  $\mathcal{NP}$ -hard problem that takes as input a graph  $G$  and a parameter  $k$ , and asks whether there exist at most  $k$  edges that can be added to the graph that would turn it into a chordal graph, i.e. one with no induced cycles of length more than 3. This problem has applications in ordering a Gaussian elimination of a matrix.

**Claim 2.4** *The minimum fill-in problem with parameter  $k$  has witness retrievable compression.*

**Proof:** Kaplan, Shamir and Tarjan [KST99] prove that this problem is fixed-parameter tractable (this notion of tractability in parameterized complexity means that the problem is polynomial-time solvable when  $k$  is sufficiently small, and in particular for all fixed  $k$ ). Their algorithm partitions the graph into two sets of nodes  $A$  and  $B$  where  $A$  is of size  $k^3$  and there are no chordless cycles (i.e. an induced cycle of length greater than 3) in  $G$  that contain vertices in  $B$ . The complexity of this partition is  $O(k^2|V||E|)$ . They then prove that  $G$  has a  $k$  edge fill-in if and only if the graph induced by  $A$  has a  $k$  edge fill-in.

Thus the compression algorithm follows the same partitioning and stores only the graph induced by the small set  $A$ . The new graph has at most  $k^3$  vertices and thus can be presented by only  $\text{poly}(k) \log |k|$  bits. The fill-in for the new instance is exactly that of the original instance and thus the compression can be witness retrievable if the original labels of the vertices are used for the compressed graph as well.  $\square$

This use of an algorithm from parameterized complexity is not a coincidence. The “problem kernel” method (see [DF99], chapter 3) is to first reduce the problem to a small sub-instance that, like compression, contains the answer to the original problem. Then the algorithm runs in exponential time algorithm on this small instance. As was discussed in Section 1.4, if the running time of the first reduction happens to be only polynomial in the parameter, then the first phase of the algorithm is a compression algorithm.

In this context, it is important to note that a compression algorithm for a language *does not* necessarily give a parameterized complexity algorithm for the same language. At first glance it seems that one can first run the compression algorithm, and then solve the compressed problem by brute force and thus get a fixed parameter algorithm. However, such a strategy does not work since in the compression algorithm the witness is allowed to grow by a factor of  $\text{polylog}(m)$ , and thus solving the compressed problem by brute force may require a super-polynomial time in  $m$ .

## 2.2 W-Reductions and Compression-Completeness

The few examples of compression that we have showed clearly indicate that the study of  $\mathcal{NP}$  problems with respect to compression gives a very different perspective than the traditional study of  $\mathcal{NP}$ . The reason is that the typical Karp-reduction between  $\mathcal{NP}$  problems does not distinguish between the length of the witness and the length of the instance. For example, when reducing SAT to the Clique problem, one builds a large graph from a CNF formula and asks whether or not it has a Clique of size  $k$ . However, in this new instance, the witness size<sup>12</sup> is a polynomial in  $m$  (the length of the SAT formula) rather than  $n$  (the number of variables in the formula). Thus, it is not clear how to use a compression algorithm for Clique to get a compression algorithm for SAT.

**W-reductions and compression-completeness:** In order to show that a compression algorithm for  $L'$  implies a compression algorithm for  $L$ , a more restricted type of reduction is needed. We call this a W-reduction and it is similar to a Karp-reduction but asks an extra property on the length of the witness.

<sup>12</sup>The witness for Clique is a choice of  $k$  vertices from the graph

**Definition 2.5 (W-Reduction)** For two  $\mathcal{NP}$  languages  $L$  and  $L'$  we say that  $L$  W-reduces to  $L'$  if there exist polynomials  $p_1$  and  $p_2$  and a polynomial time computable function  $f$  that takes an instance  $x$  for  $L$  and outputs an instance  $f(x)$  for  $L'$  such that:

1.  $f(x) \in L'$  if and only if  $x \in L$ .
2. If  $x$  is of length  $m$  with witness length  $n$ , then  $f(x)$  is of length  $p_1(n, m)$  with witness length only  $p_2(n, \log m)$ .

We first note that this reduction composes, that is:

**Claim 2.6** If  $L$  W-reduces to  $L'$  and  $L'$  W-reduces to  $L''$  then  $L$  W-reduces to  $L''$ .

We next claim that W-reduction indeed fulfils its goal with respect to compression:

**Claim 2.7** Let  $L$  and  $L'$  be  $\mathcal{NP}$  languages such that  $L'$  W-reduces to  $L$ . Then given a compression algorithm for  $L$ , one can obtain a compression algorithm for  $L'$ .

**Proof:** Suppose that  $x$  is an instance for language  $L'$  of length  $m$  with witness length  $n$ . The compression algorithm for  $L'$  runs as follows: First use the W-reduction to  $L$  and get an instance  $f(x)$  for  $L$ , and then run the compression algorithm for  $L$  on  $f(x)$ . By the properties of the reduction  $f(x)$  is of length  $m' = p_1(n, m)$  with witness length  $n' = p_2(n, \log m)$ . The outcome of the compression is therefore of length  $\text{poly}(n', \log m') = \text{poly}(n, \log m)$ . Furthermore, this outcome is in some  $\mathcal{NP}$  language  $L''$  if and only if  $f(x) \in L$  which in turn happens if and only if  $x \in L'$ . Thus the combined process gives a compression algorithm for instances of  $L'$ .  $\square$

We remark that in the complexity discussion of compression we choose to ignore the issue of witness retrievability. Nevertheless, in order for a W-reduction to relay this property as well, the reduction itself must also have a witness retrievability property. That is, given a witness  $w$  for  $x \in L$  then one can efficiently compute  $w'$  for  $f(x) \in L'$  (without the knowledge of  $x$ ). We next define complete problems with respect to compression. These are defined in standard fashion, only with respect to W-reductions.

**Definition 2.8 (Compression-Complete)** A problem  $L$  is compression-complete for class  $\mathcal{C}$  if:

1.  $L \in \mathcal{C}$
2. For every  $L' \in \mathcal{C}$  the language  $L'$  W-reduces to  $L$ .

A language is called compression-hard for class  $\mathcal{C}$  if only requirement 2 holds.

The relevance of compression-complete problem is stated in the following simple claim.

**Claim 2.9** Let  $L$  be compression-complete for class  $\mathcal{C}$ , then given a compression algorithm for  $L$ , one can obtain a compression algorithm for any  $L' \in \mathcal{C}$ .

The proof follows directly from the definition of completeness and Claim 2.7.

## 2.3 The $\mathcal{VC}$ Classification

We turn to introduce a new classification that arises from the study of compression algorithms. In order to introduce the classification we define a series of  $\mathcal{NP}$  languages. First we introduce some notation: By a **circuit of depth  $k$**  we mean a depth  $k$  alternating AND-OR circuit where the fan-in of the gates is bounded only by the size of the circuit and negations are only on the input variables (no NOT gates).

### Definition 2.10 (Depth $_k$ CircuitSAT)

For any  $k \geq 2$  consider the  $\mathcal{NP}$  problem called Depth $_k$ CircuitSAT:

**Input:** a circuit  $C$  of size  $m$  and depth at most  $k$  over  $n$  variables.

**Membership:**  $C \in \text{Depth}_k\text{CircuitSAT}$  if there exists a satisfying assignment to  $C$ .

The next language, LocalCircuitSAT, is a less natural one. It is designed to capture computations that do not need to access the whole input, but can rather check only a sub-linear fraction of the input. Let  $x$  be a string of length  $m$ , if  $I = (i_1, \dots, i_n)$  is a list of  $n$  locations in  $x$  then we denote by  $x(I)$  the values of  $x$  at these locations.

### Definition 2.11 (LocalCircuitSAT)

**Input:** A string  $x$  of length  $m$  and a circuit  $C$  over  $n + n \cdot \log m$  variables and of size  $(n + n \cdot \log m)$ .<sup>13</sup>

**Membership:** If there exists a list  $I$  of  $n$  locations in  $x$  such that  $C(x(I), I) = 1$ .

We can now introduce our classification of  $\mathcal{NP}$  problems:

**Definition 2.12 (The  $\mathcal{VC}$  classification of  $\mathcal{NP}$  problems)** Consider  $\mathcal{NP}$  problems where  $m$  denotes the instance size and  $n$  denotes the witness size. We define the class  $\mathcal{VC}_k$  for every  $k \geq 0$ . The definition is divided into three cases:

- $k = 0$ : The class  $\mathcal{VC}_0$  is the class of all languages that admit compression algorithms.
- $k = 1$ : The class  $\mathcal{VC}_1$  is the class of all languages that W-reduce to LocalCircuitSAT.
- $k \geq 2$ : The class  $\mathcal{VC}_k$  is the class of all languages that W-reduce to Depth $_k$ CircuitSAT.

For any function  $k(m, n)$  (where  $k(m, n) \leq m$ ) also define  $\mathcal{VC}_{k(m, n)}$  as the class of all languages that W-reduce to Depth $_{k(m, n)}$ CircuitSAT. Finally, define  $\mathcal{VC} = \mathcal{VC}_m$  (the class for  $k(m, n) = m$ ).

A first observation is that simply by definition, the languages LocalCircuitSAT and Depth $_k$ CircuitSAT are compression-complete for their respective classes. The most notable example is for the class  $\mathcal{VC} = \mathcal{NP}$  where the complete problem is CircuitSAT (satisfiability of a polynomial size circuit).

When discussing a W-reduction to a depth  $k$  circuit, we can actually assume without loss of generality that the top gate of this circuit is an AND gate. Formally, let Depth $_k$ CircuitSAT $_{AND}$  denote the language Depth $_k$ CircuitSAT when restricted to circuits where the top gate is an AND gate.

**Claim 2.13** For any  $k \geq 2$ , we have that a language  $L \in \mathcal{VC}_k$  if and only if  $L$  W-reduces to Depth $_k$ CircuitSAT $_{AND}$ .

<sup>13</sup>The choice of the circuit of size to be  $n'$  (over  $n'$  variables) is arbitrary and other polynomial functions suffice as well. Furthermore, such a circuit of small size may be meaningful since not all the variables have to be used and some might be just dummy variables.

**Proof:** We show that any instance that contains a circuit where the top gate is an OR W-reduces to an instance with top gate AND. We prove this first for  $k \geq 3$ . Denote the input circuit by  $C = \bigvee_j \bigwedge_t C_{j,t}$  where each  $C_{j,t}$  is a top OR depth  $(k-2)$  circuit. If  $C$  is satisfiable then  $\bigwedge_t C_{j,t}$  is satisfiable for at least one choice of  $j$ . Add to the witness the index  $i$  of this satisfiable sub-circuit ( $i$  is given by the boolean variables  $i_1, \dots, i_\ell$  where  $\ell$  is logarithmic in  $\text{poly}(m, n)$ ). For each  $j$ , denote  $C'_{j,t} = C_{j,t} \vee i_1^{\bar{j}_1} \vee \dots \vee i_\ell^{\bar{j}_\ell}$ , where  $i^{\bar{j}}$  denotes  $i \oplus j$ . Notice that  $C'_{j,t}$  is always satisfied for  $j \neq i$ , and for  $j = i$  is satisfied if and only if  $C_{j,t}$  is satisfied. Thus the circuit can now be written as  $C' = \bigwedge_{j,t} C'_{j,t}$  that is satisfiable if and only if the original circuit was. The top OR gate of  $C$  is therefore removed in the new instance  $C'$  while adding only a small number of variables, thus the input to the circuit witness remains of order  $\text{poly}(n, \log m)$  as required.

In the case  $k \geq 3$ , the depth of the new instance becomes  $k-1$ . In the case that  $k=2$ , the bottom level that included only variables is transformed into an OR of variables, thus the new circuit is simply a CNF formula (and the depth remains  $k=2$ ).  $\square$

An immediate corollary is that SAT (that is, satisfiability of CNF formulas) is compression complete for the class  $\mathcal{VC}_2$ .

**The  $\mathcal{VC}$  Hierarchy:** The  $\mathcal{VC}$  classification indeed defines a hierarchical structure. That is:

$$\mathcal{VC}_0 \subseteq \mathcal{VC}_1 \subseteq \mathcal{VC}_2 \subseteq \mathcal{VC}_3 \cdots \subseteq \mathcal{VC}.$$

And in general, for every two polynomially bounded functions  $k(n, m), \ell(n, m)$  such that for all  $n, m$  we have  $k(n, m) \leq \ell(n, m)$  then  $\mathcal{VC}_k(m, n) \subseteq \mathcal{VC}_\ell(m, n)$ . Furthermore,  $\mathcal{VC} = \mathcal{NP}$  by the definition of  $\mathcal{NP}$ . These observations follow trivially by the definitions, the only non-trivial part being the fact that  $\mathcal{VC}_1 \subseteq \mathcal{VC}_2$ , that is proved next.

**Lemma 2.14**  $\mathcal{VC}_1 \subseteq \mathcal{VC}_2$

**Proof:** We need to show a W-reduction from LocalCircuitSAT to SAT. The input is therefore a long string  $x$  and small circuit  $C$  on  $n + n \log m$  variables. Let  $i_1, \dots, i_n$  denote the potential locations in the string that the circuit  $C$  receives as inputs. Also define the variables  $y_1, \dots, y_n$  to indicate the values of  $x$  in the corresponding locations (that is  $y_t = x_{i_t}$  for  $t \in [n]$ ). Thus the circuit  $C$  runs on the variables  $y_1, \dots, y_n$  and the bits of  $i_1, \dots, i_n$ .

We first note that  $C$  is of size  $p(n, \log m) = (n + n \log m)$  and may be reduced (via Cook's Theorem [Coo71]) to a CNF formula  $\Phi_C$  over  $O(p(n, \log m))$  variables and of size  $O(p(n, \log m))$  that is satisfiable if and only if  $C$  is satisfiable.

Thus we have a CNF formula over the variables  $y_1, \dots, y_n, i_1, \dots, i_n$  and some extra variables. This formula's satisfiability will be equivalent to the membership of the LocalCircuitSAT instance if we manage to force the variables of  $y$  to take the values  $y_t = x_{i_t}$ . This is done by adding additional clauses to the CNF in the following manner: For simplicity we describe this only for  $y_1$ , where the same is repeated for every other  $y_t$  for  $t \in [n]$ . Define for each  $j \in [m]$  a formula  $\Phi_j = (y_1 = x_j) \vee (i_1 \neq j)$ . Notice that  $\Phi_{i_1} = 1$  if and only if  $y_1 = x_{i_1}$ . Denote the bits of  $i_1$  by  $i_{1,1}, \dots, i_{1,d}$  where  $d = \lceil \log m \rceil$ . An alternative way to write  $\Phi_j$  is as the following CNF (recall that  $i^{\bar{j}}$  denotes  $i \oplus j$ ):

$$\Phi_j = (y_1 \vee \bar{x}_j \vee i_{1,1}^{\bar{j}_1} \vee \dots \vee i_{1,d}^{\bar{j}_d}) \wedge (\bar{y}_1 \vee x_j \vee i_{1,1}^{\bar{j}_1} \vee \dots \vee i_{1,d}^{\bar{j}_d})$$

Finally, to force  $y_1 = x_{i_1}$  we simply take the new CNF to be  $\Phi_C \wedge \bigwedge_{j \in [m]} \Phi_j$ . The same is repeated to force  $y_t = x_{i_t}$  for all  $t \in [n]$ .  $\square$

## 2.4 The $\mathcal{VC}$ Classification and Verification with Preprocessing

We now present the  $\mathcal{VC}$  hierarchy from a different angle, that of the verification complexity of a language. This approach, though slightly more cumbersome than the definition via W-reductions, gives more intuition as to what it means to be in a class  $\mathcal{VC}_k$ . The new view defines the  $\mathcal{VC}$  hierarchy with respect to the verification algorithm for  $L$ , that is, the efficient procedure that takes a witness  $w$  for  $x \in L$  and verifies that it is indeed a true witness. We point out that the nature of verification algorithms may be very different when discussing different  $\mathcal{NP}$  problems. For example, in the  $k$ -Clique problem the verification algorithm needs to check only  $O(k^2)$  edges in the graph, and thus can read only a sub-linear part of the instance. In SAT, on the other hand, all the clauses in the formula must be checked when verifying a witness.

Simply looking at the verification algorithm of a language is not sufficient. For starters, classification according to verification does not distinguish between problems in  $\mathcal{P}$  that are trivially compressible and  $\mathcal{NP}$ -complete languages. Instead, we consider the notion of verification with preprocessing. This is the process for verifying that  $x \in L$  when given a witness, that also allows a preprocessing stage to the instance. Formally:

**Definition 2.15 (Verification with Preprocessing)** *Let  $L$  be an  $\mathcal{NP}$  language with instances of length  $m$  and witness length  $n$ . A pair of polynomial time algorithms  $(P, V)$  are called a verification with preprocessing for  $L$  if the following two step verification holds:*

1. **Preprocessing:**  $P$  gets an instance  $x$  and outputs a new instance  $P(x)$ .
2. **Verification:** There exists a polynomial  $p(\cdot, \cdot)$  such that  $x \in L$  if and only if there exists a witness  $w$  of length at most  $p(n, \log m)$  such that  $V(P(x), w) = 1$ .

Notice that when allowing for preprocessing, then all problems in  $\mathcal{P}$  have a pair  $(P, V)$  where  $P$  solves the problem and stores the answer while  $V$  simply relays this answer. Thus when considering the complexity of  $V$  in this definition, then easy problems indeed have very low complexity.

**The  $\mathcal{VC}$  Classification via Verification with Preprocessing:** An alternative and equivalent way to view the classes in the  $\mathcal{VC}$  hierarchy is based on the verification algorithm  $V$  in a verification with preprocessing pair  $(P, V)$ . The problems are divided into two families:

- The class  $\mathcal{VC}_1$  is the set of the languages that have very efficient verification (i.e.  $poly(n, \log m)$  rather than  $poly(n, m)$ ). We assume random access to the instance, thus such a verification algorithm only accesses a sub-linear fraction of the instance.
- The languages whose verification is not very efficient (run in time  $poly(n, m)$ ). This family is further classified into sub categories. The class  $\mathcal{VC}_k$  is the class of languages where the verification algorithm  $V$  has a representation as a depth  $k$  polynomial size circuit (polynomial in  $n$  and  $m$ ).

This definition is equivalent to the definition via W-reductions since the W-reduction to the complete problem can simply be viewed as the a preprocessing stage. In the other direction, every preprocessing stage is actually a W-reduction to the language defined by  $V$ .

## 2.5 Within $\mathcal{VC}_1$ - The class $\mathcal{VC}_{OR}$

Arguably, the most interesting class in the hierarchy is the bottom class  $\mathcal{VC}_1$ . It contains many natural problems such as Clique or sparse subset-sum that only test local properties of the input. Furthermore, it is presumably the easiest to find compression algorithms for. We further refine our hierarchy within the class

$\mathcal{VC}_1$  by introducing another class, the class  $\mathcal{VC}_{OR}$ . Consider the language  $OR(L)$  that take a large OR of small instances of a language  $L$ . Formally:

**Definition 2.16 (OR(L))**

Let  $L$  be an  $\mathcal{NP}$  language. Define the language  $OR(L)$  as follows:

**Input:**  $m$  instances  $x_1, \dots, x_m$  to the language  $L$ , each of size  $n$ .

**Membership:** If there exists  $i \in [m]$  such that  $x_i \in L$ .

Specifically the language **OR(CircuitSAT)** is defined as:

**Input:**  $m$  different circuits where each circuit is of size  $n$ .

**Membership:** If one of the  $m$  circuits is satisfiable.

This language is used to define the following class:

**Definition 2.17** The class  $\mathcal{VC}_{OR}$  is the class of all languages that W-reduce to  $OR(\text{CircuitSAT})$ .

We first note that in each of the  $m$  small instances, the instance length and witness length are polynomially related. So unlike the general case where we focussed only on short witness languages, when talking about  $OR(L)$ , any language  $L \in \mathcal{NP} \setminus \mathcal{P}$  is interesting. For example, the language  $OR(3\text{-SAT})$  is not trivially compressible. Moreover, it is compression-complete for  $\mathcal{VC}_{OR}$ .

**Claim 2.18** Let  $L$  be any  $\mathcal{NP}$ -complete language, then  $OR(L)$  is compression-complete for  $\mathcal{VC}_{OR}$ .

**Proof:** The W-reduction from  $OR(\text{CircuitSAT})$  to  $OR(L)$  simply runs the standard Karp reduction from  $\text{CircuitSAT}$  to  $L$  for each of the  $m$  circuits independently. The witness for each circuit was of at most  $n$  and is now of size  $p(n)$  for some polynomial  $p$ . In addition the witness contains an index of the instance of  $L$  that is satisfied, thus the total witness length is  $p(n) + \log m$ .  $\square$

For example, the problem  $OR(\text{Clique})$  that gets  $m$  small graphs (over  $n$  vertices) and asks whether at least one of the graphs has  $k$  sized clique (where  $k = O(n)$ ) is also compression-complete for  $\mathcal{VC}_{OR}$ . Moreover, we note the following claim that is relevant to our cryptographic applications (in Sections 3 and 4):

**Claim 2.19**  $\text{Clique}$  is compression-hard for  $\mathcal{VC}_{OR}$ .

**Proof:** The language  $OR(\text{Clique})$  W-reduces to  $\text{Clique}$  simply by taking one graph that is the union of all the small graphs in the  $OR(\text{Clique})$  instance. Clearly there is a clique in the union if and only if there is a clique in at least one sub-graph.  $\square$

A similar claim is true for all problems involving searching for a connected subgraph of size  $n$  in a graph of size  $m$  as long as the problem of deciding whether a graph of size  $p(n)$  contains such a subgraph is NP-Hard for some polynomial  $p(\cdot)$ . This is true, for instance, for the problem of whether there is a path of length  $n$ .<sup>14</sup> On the other hand we have that:

**Claim 2.20**  $\mathcal{VC}_{OR} \subseteq \mathcal{VC}_1$

---

<sup>14</sup>It is interesting to note that whereas the problem of finding a path of length  $n$  is fixed parameter tractable [AYZ95], Feige and Kilian [FK97] give indications that the  $\text{Clique}$  problem is hard for small  $n$  (via subexponential simulations).

**Proof:** This is best seen by W-reducing OR(Clique) to LocalCircuitSAT. Given graphs  $G_1, \dots, G_m$  for OR(Clique), generate the instance  $x = G_1, \dots, G_m$  and a circuit  $C$  that receives the locations of a clique in one of the graphs and checks whether they are indeed the edges in these locations form a clique (all belong to the same graph and are the edges induced by  $k$  vertices etc...). The size of the circuit is  $p(n, \log m)$  for some polynomial  $p$  since it checks only locations in  $x$  that belong to one graph (of size  $n$ ). Finally, add  $p(n, \log m)$  dummy variables to the circuit so that the circuit  $C$  has size becomes equal to the number of input variables (as is required in LocalCircuitSAT).  $\square$

Furthermore,  $\mathcal{VC}_0 \subseteq \mathcal{VC}_{OR}$ , since any compressible language can be W-reduced by the compression algorithm to a language with instance size  $p(n, \log m)$  and this instance can be reduced to CircuitSAT and viewed as an OR of a single small circuit and hence is in  $\mathcal{VC}_{OR}$ . Note that here too, one may need to add dummy variables to keep the Circuit quadratic in its input. Altogether we have that:

$$\mathcal{VC}_0 \subseteq \mathcal{VC}_{OR} \subseteq \mathcal{VC}_1.$$

## 2.6 The $\mathcal{VC}$ Classification and some $\mathcal{NP}$ Problems

In general, most of the  $\mathcal{VC}$  classification focuses on W-reductions to depth  $k$  circuits. The reasoning for this is that there is a certain tradeoff between depth and the number of variables. More precisely, we can reduce the depth of a verification circuit, but only at the price of adding additional variables (this is done using methods from Cook's Theorem [Coo71] and requires adding a variable for each gate in one intermediate level of the circuit). Since the number of variables is the focal point when discussing compression (as it coincides with the witness size), then depth turns out to be central in our classification.

Given our current state of knowledge, there are a few plausible views of the world. The two endpoint scenarios are that either there is compression for every language in  $\mathcal{NP}$  (as would be implied by a compression algorithm for CircuitSAT), or there is only compression for a few select problems, such as the examples in section 2.1. A third option is that there is a compression algorithm for some compression-complete problem in the hierarchy (say for  $\mathcal{VC}_k$ ), which would imply the collapse of all the classes below  $\mathcal{VC}_k$  to  $\mathcal{VC}_0$ .

We will briefly go over a few key classes in the hierarchy and a few examples of natural  $\mathcal{NP}$  problems and their classification (as we know it) within the  $\mathcal{VC}$  hierarchy:

**The class  $\mathcal{VC}_0$ :** Currently we know that this class contains all the languages in  $\mathcal{P}$ , languages that are already compressed by definition (such as 3-SAT), and the languages that we showed compression algorithms to (Vertex-cover, PRG-output and Minimum-fill-in).

**The class  $\mathcal{VC}_{OR}$ :** This class contains all languages  $OR(L)$  for an  $\mathcal{NP}$  language  $L$ . One natural example is the OR(SAT) problem which is actually a depth 3 circuit where the fan-in at the two bottom levels is bounded by  $n$  and only the top gate is allowed to be of greater fan-in. Some important languages in this class are those that need to be compressed in the cryptographic applications in Sections 3 and 4.

**The class  $\mathcal{VC}_1$ :** Since we are only interested in problems where the witness size  $n$  is much smaller than the instance size  $m$ , then many natural problems with this restriction are in  $\mathcal{VC}_1$ . For example, graph problems that ask whether a small graph can be embedded in a large graph are all in  $\mathcal{VC}_1$ . The Clique problem (with a clique of size  $n$ ), or Long-Path (a path of length  $n$  that does not hit any vertex twice) are such small graph embedding problems. Sparse Subset-Sum is another natural language in  $\mathcal{VC}_1$ . This language receives a set of  $m$  values and a target sum and asks whether there is a small subset for which the values add up exactly to the target sum.

**Dominating Set:** The problem asks, given a graph, whether there is a set of  $k$  vertices such that all the graph is in its neighbor set. Dominating set is in the class  $\mathcal{VC}_3$  as implied by the following verification: the witness is a set  $S$  and the algorithm tests that  $\forall$  vertex  $v \exists$  vertex  $u \in S$  such that  $(u, v)$  is in the graph. The  $\forall$  translates to an AND gate and the  $\exists$  translates to an OR gate. Finally, testing that an edge is in the graph requires an AND over the  $O(\log m)$  bits representing this edge. In total, this is a depth 3 circuit. Note that a straightforward verification of vertex cover will also yield a depth 3 circuit. However, while vertex cover is compressible and in  $\mathcal{VC}_0$ , for dominating set we are unaware a better method. In addition, dominating set is *compression-hard* for  $\mathcal{VC}_2$ . This is seen by a standard reduction of SAT to dominating set in which a SAT formula with  $n$  variables and  $m$  clauses is transformed into a graph with  $m + 3n$  vertices with the property that the graph has a dominating set of size  $n$  iff the SAT formula is satisfiable.<sup>15</sup>

**Weighted-SAT:** Given a CNF formula of length  $m$  the problem asks if it has a satisfying assignment of weight at most  $k$  (at most  $k$  variables are assigned the value 1). Unlike our previous discussions of SAT, the number of variables here is only bounded by  $m$  and the short witness simply consists of the list of all variables that receive the value 1 (that is, the witness is of length  $n = k \log m$ ). This problem serves as the basic complete problem for the parameterized complexity class  $W[1]$ , which is at the bottom of the W-hierarchy (see [DF99]). However, with regards to compressibility, we only know how to place it in the class  $\mathcal{VC}_4$ . This is shown by the following verification procedure (using the same logic as with Dominating-Set): For every witness (list)  $L$ , the algorithm tests that  $\forall$  clauses  $C$  either  $\exists$  a literal  $x \in C$  such that  $x \in L$  or  $\exists$  a negated literal  $\bar{x} \in C$  such that  $x \notin L$ . The verification of  $x \in L$  adds up to total depth 3 by testing that  $\exists y \in L$  such that  $x = y$  (where  $x = y$  is tested by an AND over the bits of  $x$  and  $y$ ). On the other hand, verifying that  $x \notin L$  requires total depth 4 as it runs  $\forall y \in L$  we have  $x \neq y$ . The overall depth is thus dominated by the negated variables and is thus 4.

**OR of (large) instances:** Consider the OR of CNF formulas over few variables (unlike the language  $OR(SAT)$  where the CNF formulas are considerably smaller than the fan-in of the OR gate). Such a language thus contains depth three circuits, but is actually in  $\mathcal{VC}_2$ , as implied by Claim 2.13.

**Integer Programming (IP):** An instance of integer programming consists of a list of  $m$  linear constraints on  $n$  integer variables with the goal of maximizing a linear target function over these  $n$  variables (under the list of constraints). Unlike its counterpart of linear programming, where the variables may take real values and is polynomial time solvable, integer programming is  $\mathcal{NP}$ -hard even when the variables are restricted to taking only the values 0 and 1 (one of Karp's original problems [Kar72]). Thus, the decision variant of integer programming, where the number of constraints is much larger than the number of variables, is interesting with respect to compression. First, compressing it is at least as hard as compressing SAT: given a SAT instance with  $n$  variables and  $m$  constraints it is simple to come up with a corresponding IP instance with  $2n$  variable and  $m$  constraints, i.e. IP is  $\mathcal{VC}_2$ -hard. On the other hand, a straightforward verification of a witness for this problem takes the proposed assignment for the  $n$  variables and checks if it satisfies each of the constraints. The verification of a linear constraint can be achieved in logarithmic depth (in  $n$ ), placing IP in  $\mathcal{VC}_k(n)$  for  $k(n) = \Omega(\log n)$ . We are unaware of a (significantly) better classification (of lower depth) for integer programming.

---

<sup>15</sup>In a nutshell, the reduction creates a triangle for each variable  $x_i$  of the formula. One of the nodes of the triangle is identified with the positive variable and another with its negation while the third is connected only to the other two. In addition, a vertex is created for each clause in the formula. Now, each literal is connected with all of the clauses that it appears in. The generated graph has a dominating set of size  $n$  iff the formula is satisfiable.

## 2.7 On Compression of Search Problems

So far, the  $\mathcal{NP}$  problems that we discussed were all decision problems, that is, they ask if  $x \in L$ , and are answered by YES or NO. When considering a specific  $\mathcal{NP}$  relation  $R_L$  associated with  $L$ ,<sup>16</sup> then the above decision problem has a natural search problem associated with it, which is to actually find a witness to  $x \in L$  with respect to the relation  $R_L$ . A solution to such a problem is an  $n$  bit string rather than just a single bit.

Loosely speaking, a compression algorithm for the search instance should produce a shorter output that contains enough information about some witness for the original problem.

**Definition 2.21 (Compression for search problem)** *A compression algorithm for an  $\mathcal{NP}$  search problem  $L$  (with respect to  $R_L$ ) is a pair of algorithms  $(Z, E)$  with a polynomial  $p(\cdot, \cdot)$ , where  $Z$  is a polynomial time compression algorithm and  $E$  is an unbounded extraction algorithm. Given an instance  $x$  with witness parameter  $n$  we should have that:*

1.  $Z(x)$  is of length at most  $p(n, \log m)$ .
2. If  $x \in L$  and there is a witness of length  $n$ , then  $E(Z(x)) = w$  where  $w$  is a witness to  $x \in L$  with respect to  $R_L$ .

We now show that a compression for decision problems also yields a compression for search problems, without an increase in the hierarchy.

**Theorem 2.22** *If a class  $\mathcal{VC}_k$  has a compression algorithm, then there is a compression algorithm for the search problem of a relation  $R_L$  of  $L \in \mathcal{VC}_k$ .*

The technique of the proof below also comes in handy in proving Theorem 4.4, regarding the application of the ability to compress, say SAT, to cryptanalysis in hybrid bounded storage model. In the following proof, a witness to  $x \in L$  refers to a witness according to the specific relation  $R_L$  associate with  $L$ .

**Proof:** Given an instance  $x$  to a language  $L$ , for any  $i \in [n]$ , consider the  $\mathcal{NP}$  problem  $L_i$  that asks whether there exist an  $n$  bit witness  $w$  to  $x \in L$  such that  $w_i = 1$  (the  $i^{\text{th}}$  bit of  $w$  is 1). The language  $L_i$  is also in  $\mathcal{VC}_k$  since its verification circuit is the same as the one for  $L$  with an additional AND to the variable  $w_i$  (this AND gate is incorporated into the top level AND of the circuit thus the depth remains  $k$ ).

Our first attempt is to compress the instance  $x$  for every  $i \in [n]$  with respect to the language  $L_i$  (denote such a compression by  $Z_{L_i}(x)$ ). Thus we store  $Z_{L_i}(x)$  for all  $i \in [n]$ , which amounts to  $n \cdot p(n, \log m)$  bits, for some polynomial  $p(n, \log m)$ , which is also in  $poly(n, \log m)$ . Now suppose that there is only a single witness  $w$  to  $x$ ; then one can extract  $w$  bit by bit, by solving the compressed instance of each bit. However, this idea fails when  $w$  is not the only witness, and we may inconsistent answers for the different bits.

The second attempt is to use the reduction of Valiant and Vazirani [VV86] to a unique witness. The idea is to choose a pairwise-independent hash function  $h$  that is appropriately shrinking, and add to the language the requirement that  $h(w) = 0$ . We use the following lemma:

**Lemma 2.23 ([VV86])** *Let  $L$  be an  $\mathcal{NP}$  language and for every  $x$  denote by  $W_x$  the set of all witnesses to  $x \in L$ . Let  $\ell$  be such that  $2^\ell \leq |W| \leq 2^{\ell+1}$ . Let  $\mathcal{H}_{\ell+2}$  be a family of pairwise independent hash functions with  $h : \{0, 1\}^n \rightarrow \{0, 1\}^{\ell+2}$  for all  $h \in \mathcal{H}_{\ell+2}$ . Then*

$$\Pr_{h \in \mathcal{H}_{\ell+2}}[|\{w : w \in W_x \text{ and } h(w) = 0\}| = 1] \geq \frac{1}{8}$$

<sup>16</sup>Let  $L$  be an  $\mathcal{NP}$  language with parameters  $m$  and  $n$ . A relation  $R_L$  associated with is a polynomial time function  $R_L : \{0, 1\}^m \times \{0, 1\}^n \rightarrow \{0, 1\}$ , such that for every  $x \in \{0, 1\}^m$  it holds that  $x \in L$  iff there exists a  $w \in \{0, 1\}^n$  such that  $R_L(x, w) = 1$ .

Consider the  $\mathcal{NP}$  language  $L^h$  where  $x \in L^h$  if it has a witness  $w$  for  $x \in L$  and  $h(w) = 0$ . We note that this language is also in  $\mathcal{VC}_k$  since the additional requirement can be computed efficiently over  $n$  variables (the hash is efficient) and by Cook's theorem this may be represented as a CNF formula over these variables plus only  $\text{poly}(n)$  additional variables. Thus adding the requirement of the hash does not add to the depth of the verification circuit for  $L$ .

Now, if we enumerate on all values of  $\ell$  then with probability at least  $\frac{1}{8}$ , for the correct  $\ell$  we will get that  $L^h$  has a unique witness and storing  $Z_{L_i^h}(x)$  for all  $i$  suffices to maintain the information about this witness. This can be repeated sufficiently many times (say  $O(n)$  times) so that with overwhelming probability, one of the attempts will indeed give a unique witness. However, this solution is also insufficient, since we have stored a list of  $O(n^2)$  compressed values ( $O(n)$  repetitions for each value of  $\ell \in [n]$ ) and we are guaranteed that with overwhelming probability one of them is a witness for  $x$  but we do not know which one (recall that we cannot store the original instance and thus cannot verify that a witness is correct).

Our final attempt succeeds in reducing the list of potential witnesses into a unique and true witness. This compression is as follows: Denote by  $L_{\bar{i}}$  the language that asks whether there exist an  $n$  bit witness  $w$  to  $x \in L$  such that  $w_i = 0$  (similar to  $L_i$  but with  $w_i$  negated). The compression of an instance  $x$  to the search problem  $L$  goes as follows:

For every  $\ell \in [n]$  repeat the following  $n$  times:

- Choose  $h \in_R \mathcal{H}_{\ell+2}$ .
- For all  $i \in [n]$  store  $Z_{L_i^h}(x)$  and  $Z_{L_{\bar{i}}^h}(x)$ .

The extraction procedure is as follows: For all  $\ell$  and  $h \in \mathcal{H}_{\ell+2}$ , solve all the compressed instance pairs. For every pair  $Z_{L_i^h}(x)$  and  $Z_{L_{\bar{i}}^h}(x)$ , if they both are negative or both are positive then ignore all values that are compressed with this  $h$ . Only if for all  $i$  we have exactly one of the instances being correct then output the  $i^{\text{th}}$  bit of  $w$  according to the result.

The above algorithm indeed compresses since it only adds a factor of  $n^3$  to the overall storage. With probability at least  $1 - 2^{-O(n)}$  at least one of the chosen  $h$ 's is successful in leaving exactly one witness to  $x \in L_h$ , and this witness will be extracted. Finally, if  $h$  did not leave exactly one witness, then this will be identified: If there are no witnesses then  $Z_{L_i^h}(x)$  and  $Z_{L_{\bar{i}}^h}(x)$  will both be negative for all  $i$ . If there is more than one witness, then for at least one choice of  $i$  both  $Z_{L_i^h}(x)$  and  $Z_{L_{\bar{i}}^h}(x)$  will be positive.  $\square$

**Maintaining other information:** We have seen that compression may maintain much more than just a yes/no answer. A natural question to ask what other types of information may be maintained through compression algorithms. The following are some examples:

**Number of witnesses:** The compression described above actually maintains an approximation of the number of witnesses to  $x \in L$  (with respect to  $R_L$ ). Once the chosen  $k$  is too large, there will be a sharp drop in the probability of having a witness and this can be observed when extracting the witnesses and indicate what is the right  $k$ .

**An almost random witness:** The compression above also outputs a witness that is almost uniformly distributed over  $W_x$ . Or more accurately, the probability of getting each witness is bounded by a constant times  $1/|W_x|$ .

**On maintaining all witnesses:** As opposed to maintaining a single witness or the number of witnesses, a compressed instance cannot always maintain the information about *all* of the witnesses of an input instance. This is shown by the following simple information theoretic argument: encode an  $m$  bit string  $s$  with a DNF circuit  $C$  by constructing for each position  $j \in [m]$  a formula  $C_j$  on  $\log m$

variables. If  $s[j] = 1$  then take  $C_j$  to be circuit that is satisfied iff the variables encode the index  $j$ . If  $s[j] = 0$  then  $C_j$  is the non-satisfiable circuit  $C_j = 0$ . The circuit  $C$  is formed by taking an OR of all these circuits ( $C = \bigvee_{j \in [m]} C_j$ ). The satisfying assignments of  $C$  correspond exactly to the 1's in  $s$ . Consider  $C$  as an input to the language as  $\text{CircuitSAT}$ <sup>17</sup>. Suppose that there exists a compression algorithm that maintains all of the witnesses of a circuit  $C$ . In particular, this means that the  $m$  bit string  $s$  may also be extracted from the compressed instance. But this is clearly impossible information theoretically, since  $m$  random bits may not be represented by  $\text{poly}(n, \log m) < m$  bits. So we conclude that if our goal is come up with a compression algorithm for SAT then we must come up with a way of losing information about the witnesses.

In the examples of compression that we have seen in Section 2.1, the compression algorithms for vertex cover, PRG-output and Minimum fill-in actually maintain all the witnesses. On the other hand, the compression for GapSAT (which we will see in Section 2.8) does not necessarily maintain this information, as it is based on sampling.

## 2.8 Speculation on Compression

We give two arguments that may be viewed as evidence to the existence and non-existence of compression respectively.

**An Optimistic View - Compression of a promise problem and the PCP Theorem:** Consider the promise problem GapSAT that takes as input a CNF formula  $\Phi$  of size  $m$  over  $n$  variables and the guarantee that either  $\Phi$  is satisfiable or it is at most  $(1 - \frac{1}{2n})$ -satisfiable (no assignment satisfies more than  $(1 - \frac{1}{2n})$  of its clauses). The task is to decide if  $\Phi$  is satisfiable or far from satisfiable.

Such a problem has a simple and witness retrievable compression. The idea is to choose  $O(n^2)$  random clauses from  $\Phi$  and take the AND of these clauses to be the compressed formula  $\Psi$ . This compression works because if  $\Phi$  is far from satisfiable then for every assignment the formula  $\Psi$  is satisfied with probability at most  $2^{-2n}$  ( $\Psi$  does not contain one of the  $\frac{1}{2n}m$  unsatisfied clauses). Taking a union bound over all assignments, we get that with probability  $(1 - 2^{-n})$  the formula  $\Psi$  has no satisfying assignment. On the other hand, if  $\Phi$  is satisfiable then the same assignment also satisfies  $\Psi$  (and hence the witness retrievability). Note that our definition of GapSAT is robust in the sense that GapSAT is compressible whenever the gap is  $(1 - \frac{1}{p(n)})$  for every choice of a polynomial  $p(\cdot)$ .

The above simple compression algorithm is especially interesting in light of the PCP Theorem. One way to view the PCP Theorem is as an efficient reduction from an instance of SAT to an instance of GapSAT. Thus one can hope to combine the PCP reduction with the above compression and get a compression for general SAT. However, reducing general SAT to GapSAT via the PCP is not a W-reduction as the witness size grows to the order of the instance size. For starters, the PCP Theorem is typically defined over 3-CNF formulas, and the reduction of a general size  $m$  CNF to a 3-CNF adds  $O(m)$  variables. In order for this approach to achieve compression for SAT, we require a new PCP Theorem that is actually a W-reduction.

GapSAT is just one example of a gap problem that admits compression. For instance, one can consider the promise problem GapClique where a graph of size  $m$  either has a Clique of size  $m/n$  or contains no Clique of size  $n$ . As in the case of GapSAT, GapClique is compressible by sampling a subset of its vertices. Thus, coming up with a W-reduction from a general  $(n', m')$ -Clique problem (the graph of size  $m'$  either contains a clique of size  $n'$  or not) to  $(n, m)$ -GapClique would enable the compression of Clique. We view finding PCPs that are also W-reductions as a major research direction, especially in light of the recent new proof to the PCP Theorem of Dinur [Din05].

<sup>17</sup>The circuit  $C$  is actually an instance for the language  $OR(\text{CircuitSAT})$ .

**A Pessimistic View - On Oblivious Compression:** We have shown that it is impossible to maintain all of the information in an instance when compressing it and some information is necessarily lost (for example the list of all witnesses cannot be kept). On the other hand, we show that if compression exists then it is not likely to lose too much information about the original instance. Such a result would entail the collapse of the polynomial hierarchy to its second level. More formally:

Let  $Z$  be a compression algorithm for SAT. We consider it as a two input algorithm taking a formula  $\Phi$  and local randomness  $r \in \{0, 1\}^\ell$ . Denote by  $Z(\Phi, U_\ell)$  the random variable taking the output of  $Z$  with fixed input  $\Phi$  and random  $r \in_R \{0, 1\}^\ell$ . Let  $\mathbf{X}$  be a distribution over formulas. The random variable  $Z(\mathbf{X}, U_\ell)$  denotes the output of  $Z$  under a choice of random  $r$  and a random  $\Phi$  from the distribution  $\mathbf{X}$ .

The compression  $Z$  is said to be *oblivious* if there exists a samplable distribution  $\mathbf{X}$  over satisfiable formulas, such that for every satisfiable instance  $\Phi$  the distribution  $Z(\Phi, U_\ell)$  and the distribution  $Z(\mathbf{X}, U_\ell)$  are statistically close (within statistical distance  $\varepsilon$ ).

**Claim 2.24** *If there exists an oblivious compression for SAT, then the polynomial hierarchy collapses to its second level.*

**Proof:** We show that if oblivious compression of SAT instances exists then  $\text{Co-SAT} \subseteq \mathcal{AM}$ . Consider the following interactive proof that an instance  $\Phi \notin \text{SAT}$ . The verifier chooses a random satisfiable formula  $\Psi \in \mathbf{X}$  randomness  $r \in U_\ell$  and flips a random coin  $c$ . If  $c = 0$  then the verifier sends  $\xi = Z(\Phi, r)$  to the prover, if  $c = 1$  he sends  $\xi = Z(\Psi, r)$ . The prover then answers 1 if the compressed instance is satisfiable and 0 otherwise. The verifier accepts if the provers answer equals his bit  $c$  and rejects otherwise.

*Completeness:* If indeed  $\Phi \notin \text{SAT}$ , then the prover will be able to tell whether the verifier used a coin  $c = 0$  or  $c = 1$ , simply by testing the satisfiability of  $\xi$  and replying correctly.

*Soundness:* Suppose that  $\Phi \in \text{SAT}$ , then by the obliviousness property of  $Z$  the message  $\xi$  is from nearly the same distribution whether  $c = 0$  or  $c = 1$  and the prover is bound to error with probability  $\frac{1}{2} + \varepsilon$ .  $\square$

Thus, oblivious compression for SAT is not likely to exist. However, the languages we would like to compress for the applications in Sections 3 and 4 are actually in  $\mathcal{NP} \cap \text{Co-NP}$ , and thus for these applications even oblivious compression is actually a valid possibility.

## 3 Part II: Cryptographic Applications

### 3.1 On Public Key Cryptography Based on Any One-Way Function

As mentioned in the introduction, whether one-way functions are sufficient for public key cryptography is a long standing open problem. In fact, many researchers view the black-box impossibility result of Impagliazzo and Rudich [IR89] as an indication that general one-way functions are insufficient for public key cryptography. We next show that a witness retrievable compression algorithm provides a way for bypassing this black-box impossibility.

**Theorem 3.1** *If there exists a witness retrievable compression algorithm for SAT or for Clique, or for  $\text{OR}(\text{SAT})$  (to name some examples of compression-hard for  $\mathcal{VC}_{\text{OR}}$  languages), then there exists an Oblivious Transfer (OT) protocol based on any one-way function.*

**Proof:** The construction actually builds a Private Information Retrieval (PIR) protocol, and then uses the construction of Di Crescenzo, Malkin and Ostrovsky [DMO00] to build an OT protocol from the PIR protocol. Recall that a PIR protocol has a sender with a database of size  $m$  and a receiver that chooses to learn one entry from the database. It is required that the receiver learns the bit of his choice, but a computationally

bounded sender learns essentially nothing about this choice. In addition, the total communication should be strictly smaller than  $m$ .

Suppose that  $(Z, W)$  is a witness retrievable compression algorithm for SAT (see Definition 1.2), and let  $f$  be a one-way function. Take  $(Commit, Verify)$  to be a statistically binding computationally hiding bit commitment scheme based on the one-way function  $f$ . Recall that the protocol  $Commit$  takes from the sender a string  $S$  and randomness  $r$  and after an interaction the receiver gets a commitment  $\sigma$ . The polynomial time algorithm  $Verify$  takes the commitment  $\sigma$  and a possible opening to value  $S'$  with randomness  $r'$  and verifies that  $S, r'$  are consistent with  $\sigma$ . One could take for example the commitment scheme of Naor [Nao91] based on the one-way function  $f$ .<sup>18</sup> In our proof we work under the assumption that the parties are semi-honest (that is, the parties follow the protocol as prescribed and are only allowed to try and infer extra information from the transcript of the protocol). In such a model commitment may be achieved without interaction at all. The semi-honest assumption is justified by the compiler of Goldreich, Micali and Wigderson [GMW91] that showed how to transform a semi-honest protocol into one against malicious parties (again, the only needed cryptographic assumption is the existence of a one-way function). Consider the following protocol:

**Protocol  $PIR_f$ :**

**Alice's input:** database  $D$  of  $m$  bits. Let  $D[i]$  denote the  $i$ th bit in  $D$ .

**Bob's input:** index  $i \in [m]$  denote the bits of  $i$  by  $i_1, \dots, i_\ell$

1. **Bob commits to  $i$ :** Bob commits to  $i$  with randomness  $r_B$ , Alice receives  $\sigma = Commit(i, r_B)$ .
2. **Alice computes  $\Phi$ :** The CNF formula  $\Phi$  is defined as follows:
  - Denote by  $Verify_\sigma$  the algorithm  $Verify$  with the input  $\sigma$  fixed. That is,  $Verify_\sigma$  takes as inputs  $x$  and  $r$  and accepts if and only if they form a legal opening of the commitment  $\sigma$  (and in particular this means that  $x = i$ ).
  - Translate  $Verify_\sigma$  into a CNF formula  $\Phi_\sigma$  over the variables  $x_1, \dots, x_\ell$  of  $x$  and the bits of  $r$  (using Cook's reduction).
  - For every  $j \in [m]$  define the clause  $C_j = (x_1^{j_1} \vee x_2^{j_2} \vee \dots \vee x_\ell^{j_\ell})$  if  $D[j] = 0$  (where  $x^0$  denotes  $\bar{x}$  and  $x^1$  denotes  $x$ ) and  $C_j = 1$  if  $D[j] = 1$ .
  - Set

$$\Phi = \Phi_\sigma \wedge \bigwedge_{j \in [m]} C_j$$

3. **Alice Compresses  $\Phi$ :** Alice runs  $\Psi = Z(\Phi)$  and sends  $\Psi$  to Bob.
4. **Bob checks witness:** Note that Bob knows the witness to  $Verify_\sigma$  and can compute a witness  $w$  for  $\Phi_\sigma$ . Bob checks if  $W(w, \Psi)$  is a satisfying assignment for  $\Psi$ . If it is Bob outputs 1, otherwise he outputs 0.

It remains to show that the protocol  $PIR_f$  is indeed a PIR protocol. Due to the fact that the commitment is binding (up to a negligible error), then an assignment satisfying  $\Phi_\sigma$  must have  $x = i$  (recall that  $i$  is the index that Bob committed to). Thus the first part of  $\Phi$  is only satisfied when  $x = i$ . But the second

<sup>18</sup>To be more exact, the commitment of [Nao91] can be based on the pseudorandom generator of Håstad et al. [HILL99] which in turn can be based on the function  $f$ .

part is only satisfied if  $D[x] = 1$ , thus  $\Phi$  is satisfied if and only if  $D[i] = 1$ . By the property of the compression algorithm, also  $\Psi$  is satisfiable iff  $D[i] = 1$ . Hence, using the witness retrievable properties of the compression, Bob figures out whether or not  $\Psi$  is satisfiable, and learns the bit  $D[i]$  (up to a negligible error).

The second property is that the sender Alice learns no computational information about Bob's choice. This follows directly from the guarantees of the commitment scheme (note that Bob does not send any information outside of the commitment). The third and final requirement regards the length of the communication. But the length of the communication is a fixed polynomial in  $p(n)$  (depending on the commitment protocol and the parameter of the compression algorithm). So choosing a large enough databases with  $m > p(n)$  guarantees a non trivial PIR protocol and hence an OT protocol.

Note that the proof can also go through with compression of a language in  $\mathcal{VC}_{OR}$ . This can be seen by taking for every  $j \in [m]$  such that  $D[j] = 1$  a circuit that outputs 1 if and only if there exists  $r$  such that  $Verify_\sigma(j, r)$ . Such a circuit of size polynomial in  $n$  exists due to Cook's Theorem and the OR of all these circuits equals  $D[i]$ . Thus the compression of any language that is compression-hard for  $\mathcal{VC}_{OR}$  and for which a new witness may be calculated solely based on  $Verify_\sigma$  (without involving the database) is sufficient, and in particular, a witness retrievable compression of Clique also enables the construction of OT from any one-way function (by Claim 2.19).  $\square$

We stress, as mentioned in the introduction, that compression is a valid approach towards resolving the `Minicrypt=Cryptomania` question because the protocol is inherently non-black-box, and thus does not contradict with the impossibility result of [IR89].

Note that the OT protocol derived in Theorem 3.1 is a one-round protocol (that is, one message sent from the receiver followed by one message from the sender). This follows from the construction of the PIR protocol and the construction of [DMO00] that preserves the number of rounds. One implication of this fact is that such an OT protocol may be used to construct a two round key agreement scheme, that in turn maybe used to construct a public key encryption. In general, this is achieved by fixing the first message of the protocol to be as the public key. Formally:

**Corollary 3.2** *If there exists a witness retrievable compression algorithm for SAT or for Clique, or for  $OR(SAT)$  (to name some examples of compression-hard for  $\mathcal{VC}_{OR}$  languages), then there based on any one-way function one can construct a public key encryption scheme (PKE) that is semantically secure against chosen plaintext attacks.*

### 3.2 On Collision Resistant Hash from Any One-Way Function

Loosely speaking, a collection of collision resistant hash functions (CRH) is a family  $\mathcal{H}$  of length reducing functions, such that no efficient algorithm can find collisions induced by a random hash from the family. That is, no PPTM can find for a randomly chosen  $h \in_R \mathcal{H}$ , a pair of input strings  $x$  and  $x'$  such that  $x \neq x'$  but  $h(x) = h(x')$ . In addition we want an efficient algorithm for sampling from  $\mathcal{H}$  using (possibly secret) randomness (the secret coins approach is potentially more powerful then when only public coins are used [HR04]). CRHs are important primitives with wide cryptographic applications, e.g. [Kil92, Mic94, Bar01] (see discussion and formal definitions in, for example, [IKO05]). Currently there is no known construction of CRH from general one-way functions or one-way permutations, and moreover, Simon [Sim98] shows that basing CRH on one-way permutations cannot be achieved using black-box reductions. We show, by similar techniques to the ones in the previous section, that an errorless compression of SAT (that is not necessarily witness retrievable) is sufficient to construct CRH from any one-way function.

**Theorem 3.3** *If there exists an errorless compression algorithm for SAT, or for any problem that is compression-hard for  $\mathcal{VC}_{OR}$ , then there exists a family of Collision Resistant Hash functions (CRH) based on any one-way function.*

The proof follows by combining the PIR protocol shown in the proof of Theorem 3.1, with the construction of CRH from any one-round PIR protocol due to Ishai et al. [IKO05] (note that the PIR protocol presented there is indeed a one round protocol). We give a direct proof of this theorem, with some simplified constructions and dealing with the possibility of an error with the compression.

**Proof:** As in the construction of the OT protocol in Section 3, the construction of the CRH is based on some commitment scheme which in turn may be based on any OWF [Nao91, HILL99]. A function  $h$  in the CRH collection is defined by a commitment  $\sigma$  to a value  $i \in [m]$ , and randomness  $r_Z$  for the compression algorithm. The commitment uses security parameter  $n$  where  $n \ll m$ . The hash of a string  $x$  of length  $m$  is defined as follows: For every  $j \in [m]$  let  $C_{\sigma,j}$  be the circuit that outputs one if and only if there exists randomness  $r$  such that  $\sigma$  is consistent with  $(j, r)$  (that is  $\sigma$  is a possible commitment to the value  $j$  using randomness  $r$ ). Let  $C_{\sigma,x}$  be the circuit that takes the OR of all  $C_{\sigma,j}$  such that  $x(j) = 1$  and let  $Z$  be a compression algorithm for the language  $\text{OR}(\text{CircuitSAT})$ . We define  $h_{\sigma,r_Z}(x) = Z(C_{\sigma,x}, r_Z)$ .

By the compressing properties of  $Z$  we get that  $h_{\sigma,r_Z}$  indeed shrinks its input (note that shrinkage by a single bit allows further shrinking by composition). We also have that sampling  $h_{\sigma,r_Z}$  from  $\mathcal{H}$  can be done efficiently (with secret coins).

As for collisions, let  $x \neq x'$  be two strings in  $\{0, 1\}^m$  that form a collision, i.e.,  $h_{\sigma,r_Z}(x) = h_{\sigma,r_Z}(x')$ . This equality implies that  $C_{\sigma,x}$  is satisfiable iff  $C_{\sigma,x'}$  is satisfiable. (Here we use the errorless property of the compression). We know that  $C_{\sigma,x}$  is satisfiable if and only if  $x(i) = 1$  and  $C_{\sigma,x'}$  is satisfiable if and only if  $x'(i) = 1$ . Therefore it must be the case that  $x(i) = x'(i)$ , since otherwise one of them is 0 and the other one is 1 and  $C_{\sigma,x}$  satisfiability is not that of  $C_{\sigma,x'}$ . necessarily the strings  $x$  and  $x'$  are such that  $x(i) = x'(i)$ . But for some  $j$  we have  $x(j) \neq x'(j)$  and for that  $j$  we know that  $\sigma$  is not a commitment to  $j$ .

Suppose now that we have an efficient method of finding a collision  $x$  and  $x'$  for a given  $(\sigma, r_Z)$ . Pick any  $j$  such that  $x(j) \neq x'(j)$ . Then we know that  $\sigma$  is *not* a commitment to  $j$ . This procedure can be used to break the hiding properties of the commitment scheme, since it yields an efficient method that distinguishes the commitment value from random with advantage  $1/m$ : given (the real)  $i$  and a random one  $i' \in [m]$  in a random order, run the above procedure to obtain  $j$ . If  $j$  equals one of the values in  $[m]$ , then guess this one as the random one and otherwise flip a coin. This contradicts our assumptions on building blocks.  $\square$

Note that instead of an errorless compression we can do away with an error probability slightly smaller than  $2^{-m}$ . That is, for all  $x$  we want the probability that  $Z(C_{\sigma,x}, r_Z)$  preserves the satisfiability of  $C_{\sigma,x}$  to be at least  $1 - 2^{-m+u}$  where the probability is over  $\sigma$  and  $r_Z$  and  $u \approx \log m$ . In this case we can argue (using a union bound) that with probability at least  $1 - 2^{-u}$  no  $x$  exists violating the preservation of satisfiability.

Theorem 3.3 implies the following corollary:

**Corollary 3.4** *If there exists an errorless compression algorithm for SAT or for any problem that is compression-hard for  $\mathcal{VC}_{OR}$ , then there exists statistically hiding, computationally binding commitment schemes based on any one-way function.*

The corollary follows since CRH imply statistically hiding bit commitment, see Naor and Yung [NY89] (and Damgård, Pedersen and Pfitzmann [DPP93] for commitment to many bits). As mentioned in the introduction, the current minimal assumptions for constructing statistically hiding bit commitments are the existence of one-way permutations [NOVY98] and the more general one-way functions with known pre-image size [HHK<sup>+</sup>05].

## 4 On Everlasting Security and the Hybrid Bounded Storage Model

The *bounded storage model*, introduced by Maurer [Mau92], bounds the *space* (memory size) of dishonest players rather than their running time. The model is based on a long random string  $\mathcal{R}$  of length  $m$  that is publicly transmitted and accessible to all parties. Security relies on the assumption that an adversary cannot possibly store all of the string  $\mathcal{R}$  in his memory. The requirement is that the honest parties Alice and Bob can interact using a small local storage (of size  $n$  where  $n \ll m$ ) while security is guaranteed against an eavesdropper Charlie with much larger, yet bounded storage space.

This model has enjoyed much success for the task of private key encryption. It has been shown that Alice and Bob who share a short private key can exchange messages secretly using only very small storage<sup>19</sup>, while an eavesdropper who can store up to a constant fraction of  $\mathcal{R}$  (e.g.  $\frac{1}{2}m$  bits) cannot learn anything about the messages (this was shown initially by Aumann and Rabin [AR99] and improved in [ADR02, DR02, DM04b, Lu04] and ultimately in Vadhan [Vad04]). These encryption schemes have the important property called everlasting security (put forward in [ADR02, DR02]), where once the broadcast is over and  $\mathcal{R}$  is no longer accessible then the message remains secure even if the private key is exposed and Charlie gains stronger storage capabilities.

In contrast, the situation is less desirable when Alice and Bob do not share any secret information in advance. The solution of Cachin and Maurer [CM97] for this task requires Alice and Bob to use storage of size at least  $n = \Omega(\sqrt{m})$ , which is not so appealing in this setting. Dziembowski and Maurer [DM04a] proved that this is also the best one can do.

**The Hybrid Bounded Storage Model:** The inability to achieve secure encryption in the bounded storage model with memory requirements smaller than  $n = \sqrt{m}$  has led to the following suggestion that we call the *hybrid BSM*: Let Alice and Bob agree on their secret key using a computationally secure key agreement protocol (e.g. the Diffie-Hellman protocol [DH76]). The rationale being that while an unbounded eavesdropper will eventually break the key, if this happens after the broadcast had already occurred, then the knowledge of the shared key would be useless by then (this should be expected from the everlasting security property where getting the shared key after the broadcast has ended is useless). This hybrid model is very appealing as it attempts to achieve everlasting security by adding assumptions on the ability of an adversary that has a *strict time limit*. Assumptions of this sort are generally very reasonable since all that we require is that the computational protocol is not broken in the short time period between its execution and the transmission of  $\mathcal{R}$ . For instance, an assumption such as the Diffie-Hellman key agreement [DH76] cannot be broken within half an hour, can be made with far greater degree of trust than actually assuming the long term security of this protocol.

Somewhat surprisingly, Dziembowski and Maurer [DM04a] showed that this rationale may fail. They introduce a specific computationally secure key agreement protocol (containing a non-natural modification based on private information retrieval (PIR) protocols). If this key agreement protocol is used in the hybrid BSM setting with a specific private key scheme, then the eavesdropper can completely decrypt the encrypted message. However, their result does not rule out the possibility that the hybrid idea will work with some other key agreement protocol. For instance, using the plain Diffie-Hellman key agreement may still work.

In this work we show that if compression of SAT exists then there exists an attack on the everlasting security of *any* hybrid BSM scheme.

---

<sup>19</sup>Requires  $n = O(\ell + \log m + \log \frac{1}{\epsilon})$  bits of memory for an  $\ell$  bit message and error  $\epsilon$ .

## 4.1 Two Possible Models

We define the hybrid BSM as a setting where the running time of the eavesdropper Charlie is polynomially bounded up until and during the broadcast of  $\mathcal{R}$ , and unbounded after that. We discuss two variants of a BSM scheme. We first discuss these in the standard BSM where the eavesdropper is unbounded over time, and then compare them to the hybrid setting where computational restrictions are imposed:

- **The Basic BSM Scheme:** The basic scheme does not allow interaction after the broadcast (other than sending the encrypted message). Thus the key is fully determined at the time of the broadcast. Such a scheme is fully breakable in the BSM (without an initial secret key) since the unbounded adversary can find some randomness consistent with Alice's view, and simulates Alice's actions and thus recover the encryption key<sup>20</sup>. Basic schemes in the hybrid BSM are interesting as they include any combination of a key agreement protocol with a private key scheme (such as the one described by [DM04a]). We show that if sufficiently strong compression exists then there exist attacks on any such scheme.
- **The General BSM Scheme:** Alice and Bob interact both before *and* after the broadcast. Dziembowski and Maurer [DM04a] show that such a scheme is breakable unless  $n^2 > \Omega(m)$ . For the hybrid BSM, we show that if compression exists then there exists an attack on any such scheme as long as  $n^2 > \Omega(m/p(n, \log m))$ , for some polynomial  $p$  (related to the polynomial of the compression algorithm and to the running time of the protocol that Alice and Bob use).

Thus we prove that if compression of SAT (or of any  $\mathcal{VC}_{OR}$ -hard language) is feasible then the hybrid BSM is essentially no more powerful than the standard BSM.

## 4.2 The Basic Hybrid BSM

**Definition 4.1 (Basic hybrid BSM scheme)** *A basic hybrid BSM scheme consist of the following: Alice and Bob run a protocol  $\Pi$  that is polynomial in  $n$  (this could be a key agreement scheme with security parameter  $n$ ). Denote by  $T$  the transcript of this protocol. Alice and Bob use their respective views of the protocol  $\Pi$  (i.e. the transcript  $T$  and their local randomness) to agree on  $n$  bits from the broadcast string  $\mathcal{R}$  that they should store. They store these bits and then use the stored bits to generate an encryption key  $K$  (the scheme requires that they agree on the same key).*

We show that sufficiently strong compression of SAT can be used to break any hybrid BSM scheme. For the discussion here take  $K$  to be a one bit key. The general idea is that while the eavesdropper may not figure out in time what locations to store, he can use this transcript to save a relatively short (compressed) CNF formula whose satisfiability coincides with the value of the key  $K$ . Later, when he is given unbounded computational power, he will be able to extract this bit from the compressed formula.

**Theorem 4.2** *If there exists a compression algorithm for SAT or for any compression-hard language for  $\mathcal{VC}_{OR}$ , with polynomial  $p_1$ , then any basic hybrid BSM scheme can be broken using memory  $p_2(n, \log m)$  (where  $p_2$  is a polynomial related to  $p_1$  and the running time of the protocol  $\Pi$ ).*

**Proof:** Denote the locations of the bits that Alice and Bob store by  $i_1, \dots, i_n$ . Consider the algorithm  $V$  that takes the transcript  $T_\Pi$  and the broadcast string  $\mathcal{R}$  as inputs and Alice's local randomness, and locations

---

<sup>20</sup>Since Alice must be able to decrypt the message then simulating Alice with any randomness that is consistent with the transcript must output the same key.

$i_1, \dots, i_n$  as a witness. The algorithm should check if the witness and inputs are indeed consistent with one another (for example,  $V$  should verify that a key agreement with the randomness of Alice, the transcript  $T$  indeed chooses the indices  $i_1, \dots, i_n$  to store) and output 1 if and only if they are consistent and generate an encryption key  $K = 1$ . The main observation is that the  $\mathcal{NP}$  language defined by this relation  $V$  is in  $\mathcal{VC}_1$ . Thus, if SAT has a compression algorithm then there is also a compression algorithm for all of  $\mathcal{VC}_1$  (from Lemma 2.14) including the language defined by  $V$ .

The attack of the eavesdropper Charlie is as follows: Charlie generates the verification program  $V$  and feeds the instance  $(T, \mathcal{R})$  to the compression algorithm for the language  $V$ . By the properties of the compression, the output is a CNF formula that is satisfiable if and only if  $K = 1$ . The length of the output is of some polynomial length  $p_2(n, \log m)$ . If the polynomial  $p_2$  is sufficiently small then the compressed instance is shorter than Charlie's space bound  $\frac{1}{2}m$ , and he stores this output. Finally, at a later stage, Charlie can use his unbounded powers to solve the compressed problem and retrieve the bit  $K$ .

We note that a slightly more involved argument works also with compression for  $\mathcal{V}_{\text{COR}}$ . The idea is to use independent compression for the bit  $\mathcal{R}(i_j)$  for every  $j \in [n]$ . Every such  $\mathcal{R}(i_j)$  may be presented as the OR of  $m$  circuits of size  $p(n)$  each, for some polynomial  $p$ .  $\square$

### 4.3 The General Hybrid BSM

The general scheme is like the basic one but the encryption key  $K$  is not necessarily fully defined by the end of the broadcast. In addition, the parties are allowed to interact after the broadcast is over. We note that the bounded storage key exchange scheme of Cachin and Maurer [CM97] requires such late interaction.

**Definition 4.3 (General hybrid BSM scheme)** *The general hybrid BSM scheme consist of the following: Alice and Bob run a protocol  $\Pi_1$  that is polynomial in  $n$ . Denote by  $T_1$  the transcript of this protocol. Alice and Bob use their respective views of the protocol  $\Pi_1$  to determine some  $n$  bits that each should store from the broadcast string  $\mathcal{R}$ . After the broadcast they interact in a second protocol  $\Pi_2$  (with transcript  $T_2$ ) at the end of which, both agree on encryption key  $K$ .*

**Theorem 4.4** *If there exists compression algorithm for SAT or for any compression-hard language for  $\mathcal{V}_{\text{COR}}$  with compression  $p_1(n, \log m)$ , then there exists an attack on any general hybrid BSM scheme where  $n^2 > m/p_2(n, \log m)$  (where  $p_2$  is a polynomial related to  $p_1$  and the running time of the protocol  $\Pi_1$ ).*

**Proof:** Let  $K(T_1, \mathcal{R}, T_2)$  denote the encryption key that is agreed on when the protocol is run with transcripts  $T_1, T_2$  and randomness  $\mathcal{R}$ . Because agreement is guaranteed then this key must be well defined. Denote by  $A_{T_1}$  the set of all possible randomness  $r_A$  of Alice that are consistent with the transcript  $T_1$ . Let  $s_A = S_A(T_1, \mathcal{R}, r_A)$  denote the bits that Alice stores at the end of the broadcast when running with randomness  $r_A$ , transcript  $T_1$  and broadcast string  $\mathcal{R}$ . Finally, denote by  $\mathbf{S}_A(T_1, \mathcal{R})$  the random variable that takes the value  $S_A(T_1, \mathcal{R}, r_A)$  for a uniform choice of  $r_A \in A_{T_1}$ . That is,  $\mathbf{S}_A(T_1, \mathcal{R})$  is randomly chosen from all possible  $s_A$  that Alice might have stored when running with transcript  $T_1$  and broadcast string  $\mathcal{R}$ .

We use the following important lemma of Dziembowski and Maurer [DM04a].

**Lemma 4.5 ([DM04a])** *Let  $\mathbf{S}_A(T_1, \mathcal{R})$  and  $K(T_1, \mathcal{R}, T_2)$  be defined as above. For any  $\mathcal{R}$  and  $T_1$  let  $\mathbf{S}_C(T_1, \mathcal{R})$  denote the random variables that takes  $n$  independent samples of  $\mathbf{S}_A(T_1, \mathcal{R})$ . Then:*

$$H(K(T_1, \mathcal{R}, T_2) | \mathbf{S}_C(T_1, \mathcal{R})) \leq n^2/m$$

In other words, a strategy for an eavesdropper is to store  $n$  independent samples of the random variable  $\mathbf{S}_A(T_1, \mathcal{R})$ . This strategy guarantees that the eavesdropper will have stored (with high probability) enough information on the encryption key  $K$ . Thus an eavesdropper with  $O(m)$  storage capacity may break the scheme as long as  $n^2 < O(m)$ .

Lemma 4.5 was used in [DM04a] in a setting where the eavesdropper is unbounded and can hence sample the random variable  $\mathbf{S}_A(T_1, \mathcal{R})$ . However, in our setting the eavesdropper is computationally bounded and does not have the power to generate this distribution. Instead, we use compression to store information about samples of  $\mathbf{S}_A(T_1, \mathcal{R})$  to be extracted after the broadcast is over (when the eavesdropper is unbounded).

The main idea is to use compression for search problems, as was demonstrated in Section 2.7. Define the  $\mathcal{NP}$  language  $L_A$  as follows:

$$L_A = \{(T_1, \mathcal{R}) \mid \exists \text{ witness } w = (r_A, s_A) \text{ such that } r_A \in A_{T_1} \text{ and } s_A = S_A(T_1, \mathcal{R}, r_A)\}$$

The first thing to notice is that  $L_A$  is in  $\mathcal{VC}_{OR}$ . This is shown once more by the same argument as in Theorems 4.4 or 3.1, and based on the fact that the protocol  $\Pi_1$  is polynomial time in  $n$ . Once this is established, then given a compression algorithm for  $\mathcal{VC}_{OR}$  we invoke Theorem 2.22 to get a compression algorithm to the search problem associated with  $L_A$ . Running this compression once, allows us to extract a witness to  $L_A$  and in particular to get one sample  $s_A$  of a consistent view of Alice. Running this  $n$  times supposedly gives  $n$  samples of such a view, which supposedly suffices to break the scheme by Lemma 4.5.

However, in order to invoke Lemma 4.5, we need the samples to be taken according to the distribution  $\mathbf{S}_A(T_1, \mathcal{R})$ , which is taken by a uniform distribution over  $r_A \in A_{T_1}$ . We will show that while sampling via the compression of search problems does not give the desired distribution, it is still sufficient.

A closer inspection of our compression for search technique shows that we do not necessarily sample uniformly on  $A_{T_1}$ . However, we do sample close to uniformly, in the sense that no element in  $A_{T_1}$  gets more than double the probability of another element in  $A_{T_1}$ . We then show that taking twice as many samples as was originally needed guarantees that amongst the stored bits we have  $n$  random samples of the random variable  $\mathbf{S}_A(T_1, \mathcal{R})$ , and thus we have stored enough bits from  $\mathcal{R}$  to break the scheme.

Recall from Section 2.7 that the compression algorithm for search problems chooses a random pairwise-independent hash function  $h$  and saves only a witness  $(r_A, s_A)$  that is *uniquely* hashed to the value 0 by  $h$ . Since  $r_A$  fully determines  $s_A$  (when given  $T_1$  and  $\mathcal{R}$ ) then without loss of generality we view the witness simply as  $r_A$ , furthermore, assume w.l.o.g. that  $r_A$  is of length  $n$ . Suppose that  $\ell \in [n]$  is such that  $2^\ell < |A_{T_1}| \leq 2^{\ell+1}$ . Let  $\mathcal{H}_{\ell+2}$  be a family of pairwise independent hash functions with  $h : \{0, 1\}^n \rightarrow \{0, 1\}^{\ell+2}$  for all  $h \in \mathcal{H}_{\ell+2}$ . Then for every  $r_A \in A_{T_1}$  the probability that a random  $h \in \mathcal{H}_{\ell+2}$  uniquely maps  $r_A$  to zero is at most  $2^{-(\ell+2)}$  (since  $\Pr_{h \in \mathcal{H}_{\ell+2}}[h(r_A) = 0] = 2^{-(\ell+2)}$ ). By the pairwise independence of  $\mathcal{H}$  it holds that for all other  $r'_A \in A_{T_1}$  with  $r'_A \neq r_A$  we have that  $\Pr_{h \in \mathcal{H}_{\ell+2}}[h(r'_A) \neq 0 \mid h(r_A) = 0] = 1 - 2^{-(\ell+2)}$ . By a union bound over all  $r'_A \in A_{T_1}$  with  $r'_A \neq r_A$ , combined with the probability that  $h(r_A) = 0$ , we get:

$$\Pr_{h \in \mathcal{H}_{\ell+2}}[h \text{ uniquely maps } r_A \text{ to } 0] \geq 2^{-(\ell+2)} \cdot \frac{1}{2} = 2^{-(\ell+3)}.$$

Altogether, for all  $r_A \in A_{T_1}$  it holds that

$$2^{-(\ell+2)} \geq \Pr_{h \in \mathcal{H}_{\ell+2}}[h \text{ uniquely maps } r_A \text{ to } 0] \geq 2^{-(\ell+3)}.$$

Thus whenever the hash used is indeed of length  $\ell + 2$ , the probability of sampling  $r_A \in A_{T_1}$  is almost uniform (up to a factor of 2 for each element). Since we repeat the compression for every choice of  $\ell \in [n]$  then in particular samples are stored for the correct  $\ell$ .

By Lemma 2.23 we know that at least  $\frac{1}{8}$  of the repeated compressions indeed store information about a valid witness (a sample of  $r_A \in A_{T_1}$ ). Thus, choosing, say,  $9n$  independent  $h \in \mathcal{H}_{\ell+2}$  guarantees at least  $n$  samples (by a Chernoff bound, as the choices are independent). But as mentioned above, these samples are just close to uniform over  $A_{T_1}$  rather than truly uniform. The solution is to simply run more compressions, say, for  $25n$  independent choices of  $h \in \mathcal{H}_{\ell+2}$ . This would guarantee that with overwhelming probability, at least  $3n$  samples actually are stored. We show that  $3n$  samples via the unique hashing method contain  $n$  truly uniform samples of witnesses.

This last argument follows by a hypothetical method for sampling uniformly from  $A_{T_1}$ . At a first stage,  $3n$  samples are taken using the unique hashing method. Now a diluting second stage is run: Suppose that the least likely element to be sampled gets probability  $p_{min}$ . For any element  $r_A$  that is sampled with probability  $p_{r_A}$ , keep the sample with probability  $\frac{p_{min}}{p_{r_A}}$  and delete it otherwise. Thus every element is eventually chosen with the same probability  $p_{min}$ , and since  $\frac{p_{min}}{p_{r_A}} \geq \frac{1}{2}$  then at least  $n$  samples are eventually chosen (with overwhelming probability). Note that the diluting stage is not necessarily efficiently computable, but this is taken just as a mental experiment in order to show that among the  $3n$  samples, there exist  $n$  independent samples of the random variable  $\mathbf{S}_A(T_1, \mathcal{R})$ . Thus by storing  $3n$  samples via the unique hash method, we have stored enough bits from  $\mathcal{R}$  to break the key  $K$ .  $\square$

## 5 Discussion and Open Problems

The issue of compressibility and the corresponding classification introduced in this work raise many open problems and directions. The obvious one is to come up with a compression algorithm for problem like SAT (or some  $\mathcal{VC}_{OR}$  complete or hard problem). Alternatively, show why such tasks are infeasible (see discussion in Section 2.8). We have seen compressibility of some interesting  $\mathcal{NP}$  languages and hence the question is where exactly is boundary between compressibility and incompressibility. We tend to conjecture that it is in the low levels of the  $\mathcal{VC}$  hierarchy. We view PCP amplification methods such as the recent result of Dinur [Din05] as potential leads towards achieving compression. This is since these results show a natural amplification of properties on a graph, and could potentially be combined with a simple compression of promise problems (such as the example for GapSAT in Section 2.8). The main issue is doing the PCP amplification without introducing many new variables.

Short of showing a compression for general complexity classes, it would be interesting to come up with further interesting compression algorithms as well as to obtain more hardness results. For instance, is Clique or any other embedding problem complete for  $\mathcal{VC}_1$ ? Is there a natural and simple complete problem for  $\mathcal{VC}_1$ ? Also, the  $\mathcal{VC}$  hierarchy is by no means the ultimate classification with respect to compressibility. One can hope to further refine this classification, especially within the confines of  $\mathcal{VC}$ .

Regarding cryptographic applications, one issue is how essential is the witness retrievability property required from the compression algorithm in order to get the OT from one-way function result (Theorem 3.1). In particular, is it possible that every compression algorithm for SAT can be made witness retrievable?

Since we currently do not have a general compressibility result for a significant class of languages, it is important to understand what are the implications of *incompressibility*. The application to the bounded storage model can be viewed as such a statement. Another example is the previously mentioned work of Dubrov and Ishai [DI06] regarding derandomization. In order to gain confidence in an incompressibility assumption when used in a cryptographic setting it is important to come up with an *efficiently falsifiable* assumption of this nature (see [Nao03]).

Finally we feel that we have just scratched the surface of an important topic and in the future there will be other implications of compressibility or the impossibility of compression, whether in cryptography or in other areas.

**Acknowledgements:** We thank Yuval Ishai for helpful comments and specifically for pointing out that the CRH construction does not require witness retrievability. We are also grateful to Alon Rosen and Ronen Shaltiel for their comments on the presentation.

## References

- [Aar05] S. Aaronson. NP-complete problems and physical reality. *SIGACT News*, 36(1):30–52, 2005.
- [ADR02] Y. Aumann, Y.Z. Ding, and M. O. Rabin. Everlasting security in the bounded storage model. *IEEE Transactions on Information Theory*, 48(6):1668–1680, 2002.
- [AR99] Y. Aumann and M. O. Rabin. Information theoretically secure communication in the limited storage space model. In *Advances in Cryptology – CRYPTO ’99, Lecture Notes in Computer Science*, volume 1666, pages 65–79. Springer, 1999.
- [AYZ95] N. Alon, R. Yuster, and U. Zwick. Color-coding. *Journal of the ACM*, 42(4):844–856, 1995.
- [Bar01] B. Barak. How to go beyond the black-box simulation barrier. In *42nd IEEE Symposium on Foundations of Computer Science*, pages 106–115, 2001.
- [CGH04] R. Canetti, O. Goldreich, and S. Halevi. The random oracle methodology, revisited. *Journal of the ACM*, 51(4):557–594, 2004.
- [CM97] C. Cachin and U. Maurer. Unconditional security against memory-bound adversaries. In *Advances in Cryptology – CRYPTO ’97, Lecture Notes in Computer Science*, volume 1294, pages 292–306. Springer, 1997.
- [Coo71] S. A. Cook. The complexity of theorem-proving procedures. In *3rd ACM Symposium on the Theory of Computing*, pages 151–158, 1971.
- [DF99] R. Downey and M. Fellows. **Parameterized Complexity**. Springer-Verlag, 1999.
- [DH76] W. Diffie and M.E. Hellman. New directions in cryptography. *IEEE Transaction on Information Theory*, 22:644–654, 1976.
- [DI06] B. Dubrov and Y. Ishai. On the randomness complexity of efficient sampling. To appear in STOC 2006, 2006.
- [Din05] I. Dinur. The PCP theorem by gap amplification. Electronic Colloquium on Computational Complexity (ECCC), TR05-046, 2005.
- [DLN96] C. Dwork, J. Lotspiech, and M. Naor. Digital signets: Self-enforcing protection of digital information. In *28th ACM Symposium on the Theory of Computing*, pages 489–498, 1996.
- [DM04a] S. Dziembowski and U. Maurer. On generating the initial key in the bounded-storage model. In *Advances in Cryptology – EUROCRYPT ’ 2004, Lecture Notes in Computer Science*, volume 3027, pages 126–137. Springer, 2004.
- [DM04b] S. Dziembowski and U. Maurer. Optimal randomizer efficiency in the bounded-storage model. *Journal of Cryptology*, 17(1):5–26, 2004.

- [DMO00] G. Di Crescenzo, T. Malkin, and R. Ostrovsky. Single database private information retrieval implies oblivious transfer. In *Advances in Cryptology – EUROCRYPT ’ 2000, Lecture Notes in Computer Science*, volume 1807, pages 122–138. Springer, 2000.
- [DPP93] I. Damgård, T. Pedersen, and B. Pfitzmann. On the existence of statistically hiding bit commitment schemes and fail-stop signatures. In *Advances in Cryptology - CRYPTO ’93, Lecture Notes in Computer Science*, volume 773, pages 250–265. Springer, 1993.
- [DR02] Y.Z. Ding and M.O. Rabin. Hyper-encryption and everlasting security. In *Annual Symposium on Theoretical Aspects of Computer Science (STACS)*, pages 1–26, 2002.
- [FK97] U. Feige and J. Kilian. On limited versus polynomial nondeterminism. *Chicago J. Theor. Comput. Sci.*, (1):1–20, 1997.
- [GK03] S. Goldwasser and Y. Tauman Kalai. On the (in)security of the Fiat-Shamir paradigm. In *44th IEEE Symposium on Foundations of Computer Science*, pages 102–111, 2003.
- [GLM96] J. Goldsmith, M. Levy, and M. Mundhenk. Limited nondeterminism. *SIGACT News*, 27(2):20–29, 1996.
- [GMW87] O. Goldreich, S. Micali, and A. Wigderson. How to play any mental game - a completeness theorem for protocols with honest majority. In *19th ACM Symposium on the Theory of Computing*, pages 218–229, 1987.
- [GMW91] O. Goldreich, S. Micali, and A. Wigderson. Proofs that yield nothing but their validity, or all languages in NP have zero-knowledge proof systems. *Journal of the ACM*, 38(1):691–729, 1991.
- [GV87] O. Goldreich and R. Vainish. How to solve any protocol problem - an efficiency improvement. In *Advances in Cryptology - CRYPTO ’87, Lecture Notes in Computer Science*, volume 293, pages 73–86. Springer-Verlag, 1987.
- [HHK<sup>+</sup>05] I. Haitner, O. Horvitz, J. Katz, C. Koo, R. Morselli, and R. Shaltiel. Reducing complexity assumptions for statistically-hiding commitment. In *Advances in Cryptology – EUROCRYPT ’ 2005, Lecture Notes in Computer Science*, volume 3494, pages 58–77. Springer, 2005.
- [HILL99] J. Håstad, R. Impagliazzo, L. A. Levin, and M. Luby. A pseudorandom generator from any one-way function. *SIAM Journal of Computing*, 29(4):1364–1396, 1999.
- [HN05] D. Harnik and M. Naor. On everlasting security in the hybrid bounded storage model. Manuscript, 2005.
- [HR04] C. Hsiao and L. Reyzin. Finding collisions on a public road, or do secure hash functions need secret coins? In *Advances in Cryptology – CRYPTO ’04, Lecture Notes in Computer Science*, volume 3152, pages 92–105. Springer, 2004.
- [IKO05] Y. Ishai, E. Kushilevitz, and R. Ostrovsky. Sufficient conditions for collision-resistant hashing. In *2nd Theory of Cryptography Conference – (TCC ’05)*, volume 3378 of *Lecture Notes in Computer Science*, pages 445–456, 2005.
- [Imp95] R. Impagliazzo. Hard-core distributions for somewhat hard problems. In *36th IEEE Symposium on Foundations of Computer Science*, pages 538–545, 1995.

- [IPZ98] R. Impagliazzo, R. Paturi, and F. Zane. Which problems have strongly exponential complexity? In *39th IEEE Symposium on Foundations of Computer Science*, pages 653–663, 1998.
- [IR89] R. Impagliazzo and S. Rudich. Limits on the provable consequences of one-way permutations. In *21st ACM Symposium on the Theory of Computing*, pages 44–61, 1989.
- [Kar72] R. Karp. Reducibility among combinatorial problems. In **Complexity of Computer Computations**, edited by R. Miller and J. Thatcher, New York: Plenum Press, pages 85–103, 1972.
- [KF80] C. Kintala and P. Fischer. Refining nondeterminism in relativized polynomial-time bounded computations. *SIAM Journal of Computing*, 9(1):46–53, 1980.
- [Kil88] J. Kilian. Founding cryptography on oblivious transfer. In *20th ACM Symposium on the Theory of Computing*, pages 20–31, 1988.
- [Kil92] J. Kilian. A note on efficient zero-knowledge proofs and arguments. In *24th ACM Symposium on the Theory of Computing*, pages 723–732, 1992.
- [KST99] H. Kaplan, R. Shamir, and R. Tarjan. Tractability of parameterized completion problems on chordal, strongly chordal, and proper interval graphs. *SIAM Journal of Computing*, 28(5):1906–1922, 1999.
- [Lu04] C. Lu. Encryption against space-bounded adversaries from on-line strong extractors. *Journal of Cryptology*, 17(1):27–42, 2004.
- [LV97] M. Li and P. Vitányi. **An Introduction to Kolmogorov Complexity and Its Applications**, 2nd Edition. Springer Verlag, 1997.
- [Mau92] U. Maurer. Conditionally-perfect secrecy and a provably-secure randomized cipher. *Journal of Cryptology*, 5(1):53–66, 1992.
- [Mic94] S. Micali. CS proofs. In *35th IEEE Symposium on Foundations of Computer Science*, pages 436–453, 1994.
- [MRH04] U. Maurer, R. Renner, and C. Holenstein. Indifferentiability, impossibility results on reductions, and applications to the random oracle methodology. In *The 1st Theory of Cryptography Conference – (TCC '04)*, volume 2951 of *Lecture Notes in Computer Science*, pages 21–39, 2004.
- [Nao91] M. Naor. Bit commitment using pseudorandomness. *Journal of Cryptology*, 4(2):151–158, 1991.
- [Nao03] M. Naor. On cryptographic assumptions and challenges. In *Advances in Cryptology – CRYPTO '03*, *Lecture Notes in Computer Science*, volume 2729, pages 96–109. Springer, 2003.
- [NN93] J. Naor and M. Naor. Small-bias probability spaces: Efficient constructions and applications. *SIAM Journal on Computing*, 22(4):838–856, 1993.
- [NOVY98] M. Naor, R. Ostrovsky, R. Venkatesan, and M. Yung. Perfect zero-knowledge arguments for NP using any one-way permutation. *Journal of Cryptology*, 11(2):87–108, 1998.
- [NY89] M. Naor and M. Yung. Universal one-way hash functions and their cryptographic applications. In *21st ACM Symposium on the Theory of Computing*, pages 33–43, 1989.

- [Pap94] C. Papadimitriou. **Computational Complexity**. Addison-Wesley, 1994.
- [PY88] C. Papadimitriou and M. Yannakakis. Optimization, approximation, and complexity classes. In *20th ACM Symposium on the Theory of Computing*, pages 229–234, 1988.
- [PY96] C. Papadimitriou and M. Yannakakis. On limited nondeterminism and the complexity of the V-C dimension. *Journal of Computer and System Sciences (JCSS)*, 53(2):161–170, 1996.
- [Sim98] D. Simon. Finding collisions on a one-way street: Can secure hash functions be based on general assumptions? In *Advances in Cryptology – EUROCRYPT ’ 1998, Lecture Notes in Computer Science*, volume 1403, pages 334–345. Springer, 1998.
- [TVZ04] L. Trevisan, S. Vadhan, and D. Zuckerman. Compression of samplable sources. In *IEEE Conference on Computational Complexity*, pages 1–14, 2004.
- [Vad04] S.P. Vadhan. Constructing locally computable extractors and cryptosystems in the bounded storage model. *Journal of Cryptology*, 17(1):43–77, 2004.
- [VV86] L. Valiant and V. Vazirani. NP is as easy as detecting unique solutions. *Theor. Comput. Sci.*, 47(3):85–93, 1986.
- [Wee04] H. Wee. On pseudoentropy versus compressibility. In *IEEE Conference on Computational Complexity*, pages 29–41, 2004.
- [Yao82] A. C. Yao. Protocols for secure computations. In *23rd IEEE Symposium on Foundations of Computer Science*, pages 160–164, 1982.