On the Implementation of Huge Random Objects

Oded Goldreich^{*†}

Shafi Goldwasser^{*†‡}

Asaf Nussboim^{*}

November 2, 2004

Abstract

We initiate a general study of pseudo-random implementations of huge random objects, and apply it to a few areas in which random objects occur naturally. For example, a random object being considered may be a random connected graph, a random bounded-degree graph, or a random error-correcting code with good distance. A pseudo-random implementation of such type T objects must generate objects of type T that can not be distinguished from random ones, rather than objects that can not be distinguished from type T objects (although they are not type T at all).

We will model a type T object as a function, and access objects by queries into these functions. We investigate supporting both standard queries that only evaluates the primary function at locations of the user's choice (e.g., edge queries in a graph), and complex queries that may ask for the result of a computation on the primary function, where this computation may be infeasible to perform with a polynomial number of standard queries (e.g., providing the next vertex along a Hamiltonian path in the graph).

Keywords: Pseudorandomness, Random Graphs, Random Codes, Random Functions, monotone graph properties, random walks on regular graphs.

^{*}Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot, ISRAEL. Email: {oded, shafi, asafn}@wisdom.weizmann.ac.il.

[†]Supported by the MINERVA Foundation, Germany.

[‡]Laboratory for Computer Science, MIT.

Contents

1	Introduction 2				
2	Formal Setting and General Observations2.1Specification	$ \begin{array}{c} 4 \\ 4 \\ 5 \\ 8 \\ 9 \\ 11 \end{array} $			
3	Our Main Results 3.1 Truthful Implementations . 3.1.1 Supporting complex queries regarding boolean functions . 3.1.2 Supporting complex queries regarding length-preserving functions . 3.1.3 Random graphs of various types . 3.1.4 Supporting complex queries regarding random graphs . 3.1.5 Random bounded-degree graphs of various types . 3.2 Almost-Truthful Implementations . 3.2.1 Random codes of large distance . 3.2.2 Random graphs of various types .	 13 13 14 14 15 16 16 17 17 			
4	Implementing Random Codes of Large Distance 17				
5	Boolean Functions and Interval-Sum Queries 19				
6	Random Graphs Satisfying Global Properties 2 6.1 Truthful implementations 2 6.2 Almost-truthful implementations 2				
7	Supporting Complex Queries regarding Random Graphs 31				
8	Random Bounded-Degree Graphs and Global Properties	34			
9	Supporting Complex Queries regarding Length-Preserving Functions 39				
10	Conclusions and Open Problems	43			
Bi	ibliography	45			
Aŗ	ppendix A: Implementing various probability distributionsA.1Sampling the binomial distributionA.2Sampling from the two-set total-sum distributionA.3A general tool for sampling strange distributions	47 47 48 49			
Aŗ	Appendix B: Implementing a Random Bipartite Graph 5				
Aŗ	Appendix C: Various Calculations				
Ap	Appendix D: A strengthening of Proposition 2.15				

1 Introduction

Suppose that you want to run some experiments on random codes (i.e., subsets of $\{0,1\}^n$ that contain $K = 2^{\Omega(n)}$ strings). You actually take it for granted that the random code will have *large* (i.e., linear) *distance*, because you know some Coding Theory and are willing to discard the negligible probability that a random code will not have a large distance. Suppose that you want to be able to keep succinct representations of these huge codes and/or that you want to generate them using few random bits. Being aware of the relevant works on pseudorandomness (e.g., [20, 5, 33, 16]), you plan to use pseudorandom functions [16] in order to efficiently generate and store representations of these codes; that is, using the pseudorandom function $f : [K] \to \{0, 1\}^n$, you can define the code $C_f = \{f(i) : i \in [K]\}$, and efficiently produce codewords of C_f . But wait a minute, do the codes that you generate this way have a large distance?

The point is that having a large distance is a global property of the code, which in turn is a huge (i.e., $\exp(n)$ -sized) object. This global property cannot be decided by looking at polynomially many (i.e., poly(n)-many) codewords, and so its violation cannot be translated to a contradiction of the pseudorandomness of the function. Indeed, the substitution of a random function (or a random code) by a pseudorandom one is not guaranteed to preserve the global property. Specifically, all pseudorandom codes generated as suggested above may have small distance.¹

So, can we efficiently generate random-looking codes of large distance? Specifically, can we provide a probabilistic polynomial-time procedure that allows to sample codewords from a code of large distance such that the sampled codewords look as if they were taken from a random code (which, in particular, means that we do not generate linear codes). The answer is essentially positive: see Section 4. However, this is merely an example of the type of questions that we deal with. (Another illustrative example is presented in Section 6.)

We initiate a general study of the feasibility of implementing (huge) random objects. For a given Type T of objects, we aim at generating pseudorandom objects of Type T. That is, we want the generated object to always be of Type T, but we are willing to settle for Type T objects that look as if they are truly random Type T objects (although they are not). We stress that our focus is on Type T objects that look like random Type T objects, rather than objects that look like random Type T objects although they are not of Type T at all. For example, we disapprove of a random function as being an implementation of a random permutation, although the two look alike to anybody restricted to resources that are polynomially related to the length of the inputs to the function. Beyond the intuitive conceptual reason for the above disapproval, there are practical considerations. For example, if somebody supplies an element in the range then we may want to be guaranteed that this element has a unique preimage (as would be the case with any permutation but not with a random function).

In general, when one deals (or experiments) with an object that is supposed to be of Type T, one may assume that this object has all the properties enjoyed by all Type T objects. If this assumption does not hold (even if one cannot detect this fact during initial experimentation) then an application that depends on this assumption may fail. One reason for the failure of the application may be that it uses significantly more resources than those used in the initial experiments that failed to detect the problem. Another issue is that the probability that the application fails may indeed be negligible (as is the probability of detecting the failure in the initial experiments), but due to the

¹Indeed, for each function f_s taken from some pseudorandom ensemble $\{f_s : [2^{|s|/10}] \to \{0,1\}^{|s|}\}_s$, it may hold that the Hamming distance between $f_s(i_s)$ and $f_s(i_s + 1)$ is one, for some i_s that depends arbitrarily on f_s . For example, given a pseudorandom ensemble $\{f_s\}$, consider the ensemble $\{f_{s,i}\}$ such that $f_{s,i}(i) = 0^n$, $f_{s,i}(i+1) = 0^{n-1}1$ and $f_{s,i}(x) = f_s(x)$ for all other x's.

importance of the application we are unwilling to tolerate even a negligible probability of failure.

We explore several areas in which the study of random objects occurs naturally. These areas include graph theory, coding theory and cryptography. We provide implementations of various natural random objects, which were considered before in these areas (e.g., the study of random graphs [6]).

Objects, specifications, implementations and their quality

Our focus is on huge objects; that is, objects that are of size that is exponential in the running time of the applications. Thus, these (possibly randomized) applications may inspect only small portions of the object (in each randomized execution). The object may be viewed as a function (or an oracle), and inspecting a small portion of it is viewed as receiving answers to a small number of adequate queries. For example, when we talk of huge dense graphs, we consider adjacency queries that are vertex-pairs with answers indicating whether or not the queried pair is connected by an edge. When we talk of huge bounded-degree graphs, we consider incidence queries that correspond to vertices with answers listing the neighbors of the queried vertex.

We are interested in classes of objects (or object types), which can be viewed as classes of functions. (Indeed, we are not interested in the trivial case of generic objects, which is captured by the class of all functions.) For example, when we talk of simple undirected graphs in the adjacency predicate representation, we only allow symmetric and non-reflexive Boolean functions. Similarly, when we talk of such bounded-degree graphs in the incident-lists representation, we restrict the class of functions in a less trivial manner (i.e., u should appear in the neighbor-list of v iff v appears in the neighbor-list of u). More interestingly, we may talk of the class of connected (or Hamiltonian) graphs, in which case the class of functions is even more complex. This formalism allows to talk about objects of certain types (or of objects satisfying certain properties). In addition, it captures complex objects that support "compound queries" to more basic objects. For example, we may consider an object that answers queries regarding a global property of a Boolean function (e.g., the parity of all the function's values). The queries may also refer to a large number of values of the function (e.g., the parity of all values assigned to arguments in an interval that is specified by the query).

We study probability distributions over classes of objects. Such a distribution is called a specification. Formally, a specification is presented by a *computationally-unbounded* probabilistic Turing machine, where each setting of the machine's random-tape yields a huge object. The latter object is defined as the corresponding input-output relation, and so queries to the object are associated with inputs to the machine. We consider the distribution on functions obtained by selecting the specification's random-tape uniformly. For example, a random N-vertex Hamiltonian graph is specified by a computationally-unbounded probabilistic machine that uses its random-tape to determine such a (random Hamiltonian) graph, and answers adjacency queries accordingly. Another specification may require to answer, in addition to adjacency queries regarding a uniformly selected N-vertex graph, also more complex queries such as providing a clique of size $\log_2 N$ that contains the queried vertex. We stress that the specification is not required to be even remotely efficient (but for sake of simplicity we assume that it is recursive).

Our ultimate goal will be to provide a *probabilistic polynomial-time* machine that implements the desired specification. That is, we consider the *probability distribution* on functions induced by fixing of the random-tape of the latter machine in all possible ways. Again, each possible fixing of the random-tape yields a function corresponding to the input-output relation (of the machine per this contents of its random-tape).

Indeed, a key question is how good is the implementation provided by some machine. We consider two aspects of this question. The first (and more standard) aspect is whether one can distinguish the implementation from the specification when given oracle access to one of them. Variants include perfect indistinguishability, statistical-indistinguishability and computational-indistinguishability. We stress a second aspect regarding the quality of implementation: the truthfulness of the implementation with respect to the specification, where being truthful means that any possible function that appears with non-zero probability in the implementation must also appear with non-zero probability in the specification. For example, if the specification is of a random Hamiltonian graph then a truthful implementation must always yield a Hamiltonian graph. (A reasonable relaxation of the notion of truthfulness is to require that all but a negligible part of the probability mass of the implementation is assigned to functions that appear with non-zero probability in the specification; an implementation satisfying this relaxation is called **almost-truthful**.)

Truthfulness and beyond

Our presentation highlights the notion of truthfulness, and we have tried to justify the importance we attach to this notion. Nevertheless, we wish to stress that this paper also initiates the study of general implementations, regardless of truthfulness. That is, we ask which specifications have implementations (which are indistinguishable from them). We also stress that some of our constructions are interesting regardless of their truthfulness.

Organization

In Section 2, we present formal definitions of the notions discussed above as well as basic observations regarding these notions. These are followed by a few known examples of non-trivial implementations of various random objects (which are retrospectively cast nicely in our formulation). In Section 3, we state a fair number of new implementations of various random objects, while deferring the constructions (and proofs) to the corresponding sections. These implementations demonstrate the applicability of our notions to various domains such as functions, graphs and codes. Conclusions and open problems are presented in Section 10.

2 Formal Setting and General Observations

Throughout this work we let n denote the feasibility parameter. Specifically, feasible-sized objects have an explicit description of length poly(n), whereas huge objects have (explicit description) size exponential in n. The latter are described by functions from poly(n)-bit strings to poly(n)-bit strings. Whenever we talk of efficient procedures we mean algorithms running in poly(n)-time.

2.1 Specification

A huge random object is specified by a *computationally-unbounded* probabilistic Turing machine. For a fixed contents of the random-tape, such a machine defines a (possibly partial) function on the set of all binary strings. Such a function is called an instance of the specification. We consider the input-output relation of this machine when the random-tape is uniformly distributed. Loosely speaking, this is the random object specified by the machine.

For sake of simplicity, we confine our attention to machines that halt with probability 1 on every input. Furthermore, we will consider the input-output relation of such machines only on inputs of some specified length ℓ , where ℓ is always polynomially related to the feasibility parameter n.

Thus, for such a probabilistic machine M and length parameter $\ell = \ell(n)$, with probability 1 over the choice of the random-tape for M, machine M halts on every $\ell(n)$ -bit long input.

Definition 2.1 (specification): For a fixed function $\ell: \mathbb{N} \to \mathbb{N}$, the instance specified by a probabilistic machine M, random-tape ω and parameter n is the function $M_{n,\omega}$ defined by letting $M_{n,\omega}(x)$ be the output of M on input $x \in \{0,1\}^{\ell(n)}$ when using the random-tape $\omega \in \{0,1\}^{\infty}$. The random object specified by M and n is defined as $M_{n,\omega}$ for a uniformly selected $\omega \in \{0,1\}^{\infty}$.

Note that, with probability 1 over the choice of the random-tape, the random object (specified by M and n) depends only on a finite prefix of the random-tape. Let us clarify our formalism by casting in it several simple examples, which were considered before (cf. [16, 29]).

Example 2.2 (a random function): A random function from n-bit strings to n-bit strings is specified by the machine M that, on input $x \in \{0,1\}^n$ (parameter n and random-tape ω), returns the $\operatorname{idx}_n(x)$ -th n-bit block of ω , where $\operatorname{idx}_n(x)$ is the index of x within the set of n-bit long strings.

Example 2.3 (a random permutation): Let $N = 2^n$. A random permutation over $\{0,1\}^n \equiv [N]$ can be specified by uniformly selecting an integer $i \in [N!]$; that is, the machine uses its randomtape to determine $i \in [N!]$, and uses the *i*-th permutation according to some standard order. An alternative specification, which is easier to state (alas even more inefficient), is obtained by a machine that repeatedly inspect the N next n-bit strings on its random-tape, until encountering a run of N different values, using these as the permutation. Either way, once a permutation π over $\{0,1\}^n$ is determined, the machine answers the input $x \in \{0,1\}^n$ with the output $\pi(x)$.

Example 2.4 (a random permutation coupled with its inverse): In continuation to Example 2.3, we may consider a machine that selects π as before, and responds to input (σ, x) with $\pi(x)$ if $\sigma = 1$ and with $\pi^{-1}(x)$ otherwise. That is, the object specified here provides access to a random permutation as well as to its inverse.

2.2 Implementations

Definition 2.1 places no restrictions on the complexity of the specification. Our aim, however, is to implement such specifications *efficiently*. We consider several types of implementations, where in all cases we aim at *efficient* implementations (i.e., machines that respond to each possible input within polynomial-time). Specifically, we consider two parameters:

- 1. The type of model used in the implementation. We will use either a polynomial-time *oracle machine having access to a random oracle* or a standard probabilistic polynomial-time machine (viewed as a deterministic *machine having access to a finite random-tape*).
- 2. The similarity of the implementation to the specification; that is, is the implementation may be *perfect*, *statistically indistinguishable* or only *computationally indistinguishable* from the specification (by probabilistic polynomial-time oracle machines that try to distinguish the implementation from the specification by querying it at inputs of their choice).

Our real goal is to derive implementations by *ordinary machines* (having as good a quality as possible). We thus view implementations by *oracle machines having access to a random oracle* as merely a clean abstraction, which is useful in many cases (as indicated by Theorem 2.9 below).

Definition 2.5 (implementation by oracle machines): For a fixed function $\ell : \mathbb{N} \to \mathbb{N}$, a (deterministic) polynomial-time oracle machine M and oracle f, the instance implemented by M^f and parameter n is the function M^f defined by letting $M^f(x)$ be the output of M on input $x \in \{0,1\}^{\ell(n)}$ when using the oracle f. The random object implemented by M with parameter n is defined as M^f for a uniformly distributed $f : \{0,1\}^* \to \{0,1\}$.

In fact, $M^{f}(x)$ depends only on the value of f on inputs of length bounded by a polynomial in |x|. Similarly, an ordinary probabilistic polynomial-time (as in the following definition) only uses a poly(|x|)-bit long random-tape when invoked on input x. Thus, for feasibility parameter n, the machine handles $\ell(n)$ -bit long inputs using a random-tape of length $\rho(n) = \text{poly}(\ell(n)) = \text{poly}(n)$, where (w.l.o.g.) ρ is 1-1.

Definition 2.6 (implementation by ordinary machines): For fixed functions $\ell, \rho: \mathbb{N} \to \mathbb{N}$, an ordinary polynomial-time machine M and a string r, the instance implemented by M and random-tape r is the function M_r defined by letting $M_r(x)$ be the output of M on input $x \in \{0, 1\}^{\ell(\rho^{-1}(|r|))}$ when using the random-tape r. The random object implemented by M with parameter n is defined as M_r for a uniformly distributed $r \in \{0, 1\}^{\rho(n)}$.

We stress that an instance of the implementation is fully determined by the machine M and the random-tape r (i.e., we disallow "implementations" that construct the object on-the-fly while depending and keeping track of all previous queries and answers).

For a machine M (either a specification or an implementation) we identify the pair (M, n) with the random object specified (or implemented) by machine M and feasibility parameter n.

Definition 2.7 (indistinguishability of the implementation from the specification): Let S be a specification and I be an implementation, both with respect to the length function $\ell: \mathbb{N} \to \mathbb{N}$. We say that I perfectly implements S if, for every n, the random object (I, n) is distributed identically to the random object (S, n). We say that I closely-implements S if, for every oracle machine M that on input 1^n makes at most polynomially-many queries all of length $\ell(n)$, the following difference is negligible² as a function of n

$$|\mathbf{Pr}[M^{(I,n)}(1^n) = 1] - \mathbf{Pr}[M^{(S,n)}(1^n) = 1]|$$
(1)

We say that I pseudo-implements S if Eq. (1) holds for every probabilistic polynomial-time oracle machine M that makes only queries of length equal to $\ell(n)$.

We stress that the notion of a close-implementation does not say that the objects (i.e., (I, n) and (S, n)) are statistically close; it merely says that they cannot be distinguished by a (computationally unbounded) machine that asks polynomially many queries. Indeed, the notion of pseudo-implementation refers to the notion of computational indistinguishability (cf. [20, 33]) as applied to functions (see [16]). Clearly, any perfect implementation is a close-implementation, and any close-implementation is a pseudo-implementation. Intuitively, the oracle machine M, which is sometimes called a (potential) distinguisher, represents a user that employs (or experiments with) the implementation. It is required that such a user cannot distinguish the implementation from the specification, provided that the user is limited in its access to the implementation or even in its computational resources (i.e., time).

²A function $\mu : \mathbb{N} \to [0, 1]$ is called negligible if for every positive polynomial p and all sufficiently large n's it holds that $\mu(n) < 1/p(n)$.

Indeed, it is trivial to perfectly implement a random function (i.e., the specification given in Example 2.2) by using an *oracle machine* (with access to a random oracle). In contrast, the main result of Goldreich, Goldwasser and Micali [16] can be cast by saying that there exist a pseudo-implementation of a random function by an *ordinary machine*, provided that pseudorandom generators (or, equivalently, one-way function [21]) do exist. In fact, under the same assumption, it is easy to show that every specification having a pseudo-implementation by an oracle machine also has a pseudo-implementation by an ordinary machine. A stronger statement will be proven below (see Theorem 2.9).

Truthful implementations. An important notion regarding (non-perfect) implementations refers to the question of whether or not they satisfy properties that are enjoyed by the corresponding specification. Put in other words, the question is whether *each instance of the implementation is also an instance of the specification*. Whenever this condition holds, we call the implementation truthful. Indeed, every perfect implementation is truthful, but this is not necessarily the case for close-implementations. For example, a random function is a close-implementation of a random permutation (because it is unlikely to find a collision among polynomially-many preimages); however, a random function is *not* a *truthful* implementation of a random permutation.

Definition 2.8 (truthful implementations): Let S be a specification and I be an implementation. We say that I is truthful to S if for every n the support of the random object (I, n) is a subset of the support of the random object (S, n).

Much of this work is focused on truthful implementations. The following simple result is useful in the study of the latter. We comment that this result is typically applied to (truthful) *close*-implementations by *oracle* machines, yielding (truthful) *pseudo*-implementations by *ordinary* machines.

Theorem 2.9 Suppose that one-way functions exist. Then any specification that has a pseudoimplementation by an oracle machine (having access to a random oracle) also has a pseudoimplementation by an ordinary machine. Furthermore, if the former implementation is truthful then so is the latter.

The sufficient condition is also necessary, because the existence of pseudorandom functions (i.e., a pseudo-implementation of a random function) implies the existence of one-way functions. In view of Theorem 2.9, whenever we seek truthful implementations (or, alternatively, whenever we do not care about truthfulness at all), we may focus on implementations by oracle machines.

Proof: First we replace the random oracle used by the former implementation by a pseudorandom oracle (available by the results of [16, 21]). No probabilistic polynomial-time distinguisher can detect the difference, except with negligible probability. Furthermore, the support of the pseudorandom oracle is a subset of the support of the random oracle, and so the truthful property is inherited by the latter implementation. Finally, we use an ordinary machine to emulate the oracle machine that has access to a pseudorandom oracle.

Almost-Truthful implementations. Truthful implementations guarantee that each instance of the implementation is also an instance of the specification (and is thus "consistent with the specification"). A meaningful relaxation of this guarantee refers to the case that almost all the probability mass of the implementation is assigned to instances that are consistent with the specification (i.e., are in the support of the latter). Specifically, we refer to the following definition.

Definition 2.10 (almost-truthful implementations): Let S be a specification and I be an implementation. We say that I is almost-truthful to S if the probability that (I, n) is not in the support of the random object (S, n) is bounded by a negligible function in n.

Interestingly, almost-truthfulness is not preserved by the construction used in the proof of Theorem 2.9. In fact, there exists specifications that have almost-truthful close-implementations by oracle machines but not by ordinary machines (see Theorem 2.11 below). Thus, when studying almost-truthful implementations, one needs to deal directly with ordinary implementations (rather than focus on implementations by oracle-machines). Indeed, we will present a few examples of almost-truthful implementations that are not truthful.

Theorem 2.11 There exists a specification that has an almost-truthful close-implementation by an oracle machine but has no almost-truthful implementation by an ordinary machine.

We stress that the theorem holds regardless of whether or not the latter (almost-truthful) implementation is indistinguishable from the specification.

Proof: Consider the specification of a uniformly selected function $f : \{0,1\}^n \to \{0,1\}$ having (time-bounded) Kolmogorov Complexity³ greater than 2^{n-1} . That is, the specification machine scans its random-tape, looking for a block of 2^n bits of (time-bounded) Kolmogorov Complexity greater than 2^{n-1} , and once found uses this block as a truth-table of the desired Boolean function. Since all but a negligible fraction of the functions have Kolmogorov Complexity greater than 2^{n-1} , a almost-truthful close-implementation by an oracle machine may just use a random function. On the other hand, any implementation by an ordinary machine (of randomness complexity ρ) induces a function $f : \{0,1\}^n \to \{0,1\}$ of (time-bounded) Kolmogorov Complexity at most $(O(1) + \rho(n)) + \log_2(\operatorname{poly}(n) \cdot 2^n) = \operatorname{poly}(n)$. Thus, any such implementation yields a function that violates the specification, and so cannot be even "remotely" truthful.

2.3 Known non-trivial implementations

In view of Theorem 2.9, when studying truthful implementations, we focus on implementations by oracle machines. In these cases, we shorthand the phrase implementation by an oracle machine by the term implementation. Using the notion of truthfulness, we can cast the non-trivial implementation of a random permutation provided by Luby and Rackoff [29] as follows.

Theorem 2.12 [29]: There exists a truthful close-implementation of the specification provided in Example 2.3. That is, there exists a truthful close-implementation of the specification that uniformly selects a permutation π over $\{0,1\}^n$ and responses to the query $x \in \{0,1\}^n$ with the value $\pi(x)$.

Contrast Theorem 2.12 with the trivial non-truthful implementation (by a random function) mentioned above. Note that, even when ignoring the issue of truthfulness, it is non-trivial to provide a close-implementation of Example 2.4 (i.e., a random permutation along with its inverse).⁴ However, Luby and Rackoff [29] have also provided a truthful close-implementation of Example 2.4.

³Loosely speaking, the (standard) Kolmogorov Complexity of a string s is the minimum length of a program II that produce s. The time-bounded Kolmogorov Complexity of a string s is the minimum, taken over programs II that produce s, of $|\Pi| + \log_2(\texttt{time}(\Pi))$, where $\texttt{time}(\Pi)$ is the running-time of Π . We use time-bounded Kolmogorov Complexity in order to allow for a recursive specification.

⁴A random function will fail here, because the distinguisher may distinguish it from a random permutation by asking for the inverse of a random image.

Theorem 2.13 [29]: There exists a truthful close-implementation of the specification that uniformly selects a permutation π over $\{0,1\}^n$ and responses to the query $(\sigma, x) \in \{-1,+1\} \times \{0,1\}^n$ with the value $\pi^{\sigma}(x)$.

Another known result that has the flavor of the questions that we explore was obtained by Naor and Reingold [31]. Loosely speaking, they provided a truthful close-implementation of a permutation selected uniformly among all permutations having a certain cycle-structure.

Theorem 2.14 [31]: For any $N = 2^n$, t = poly(n), and $C = \{(c_i, m_i) : i = 1, ..., t\}$ such that $\sum_{i=1}^{t} m_i c_i = N$, there exists a truthful close-implementation of a uniformly distributed permutation that has m_i cycles of size c_i , for i = 1, ..., t.⁵ Furthermore, the implementation instance that uses the permutation π can also support queries of the form (x, i) to be answered by $\pi^i(x)$, for any $x \in \{0, 1\}^n$ and any integer i (which is presented in binary).

We stress that the latter queries are served in time poly(n) also in case $i \gg poly(n)$.

2.4 A few general observations

Theorem 2.11 asserts the existence of specifications that cannot be implemented in an *almost-truthful* manner by an *ordinary* machine, regardless of the level of indistinguishability (of the implementation from the specification). We can get negative results that refer also to implementations by *oracle* machines, regardless of truthfulness, by requiring the implementation to be sufficiently indistinguishable (from the specification). Specifically:

Proposition 2.15 The following refers to implementations by oracle machines and disregard the issue of truthfulness.

- 1. There exist specifications that cannot be closely-implemented.
- 2. Assuming the existence of one-way functions, there exist specifications that cannot be pseudoimplemented.

The hypothesis in Part 2 can be relaxed: It suffices to assume the existence of NP-sets for which it is feasible to generate hard instances. For details see Appendix D.

Proof: For starters, note that the specification may just disregard the issue of randomness and invert a one-way function at images of the user's choice. Certainly, this specification cannot be pseudo-implemented, because such an implementation would yield an algorithm that violates the hypothesis (of Part 2).⁶ We may easily adapt this example such that the specification gives rise to a random object. For example, the specification may state that, given a pair of strings, one should use a random function to select one of these strings, and answer with its inverse under the one-way function. A pseudo-implementation of this specification can also be shown to contradict the hypothesis. The above refers to Part 2. Turning to Part 1, we may use a function constructed in exponential-time that cannot be inverted, except for with negligible probability, by any polynomial-time machine that uses a random oracle. That is, the specification determines such a function, and

⁵Special cases include involutions (i.e., permutations in which all cycles have length 2), and permutations consisting of a single cycle (of length N). These cases are cast by $C = \{(2, N/2)\}$ and $C = \{(N, 1)\}$, respectively.

⁶Consider the performance of the specification (resp., implementation) when queried on a randomly generated image, and note that the correctness of the answer can be efficiently verified. Thus, while the specification always inverts the one-way function on the given image, the implementation must fail except with negligible probability.

inverts it at inputs of the user's choice. Observe that a close-implementation of such a function is required to successfully invert the function at random inputs, which is impossible (except for negligible probability).

The randomness complexity of implementations: Looking at the proof of Theorem 2.9, it is evident that as far as pseudo-implementations by ordinary machines are concerned (and assuming the existence of one-way functions), randomness can be reduced to any power of the feasibility parameter (i.e., to n^{ϵ} for every $\epsilon > 0$). The same holds with respect to truthful pseudo-implementations. On the other hand, the proof of Theorem 2.11 suggests that this collapse in the randomness complexity cannot occur with respect to almost-truthful implementations by ordinary machines (regardless of the level of indistinguishability of the implementation from the specification).

Theorem 2.16 (a randomness hierarchy): For every polynomial ρ , there exists a specification that has an almost-truthful close-implementation by an ordinary machine that uses a random-tape of length $\rho(n)$, but has no almost-truthful implementation by an ordinary machine that uses a randomtape of length $\rho(n) - \omega(\log n)$.

Proof: Let $g(n) = \omega(\log n)$. Consider the specification that selects uniformly a string $r \in \{0, 1\}^{\rho(n)}$ of (time-bounded) Kolmogorov Complexity at least $\rho(n) - g(n)$, and responds to the query $i \in [2^n]$ with the $(1 + (i \mod \rho(n)))$ -th bit of r. Since all but a $\exp(-g(n)) = n^{-\omega(1)}$ fraction of the $\rho(n)$ -bit long string have such complexity, this specification is closely-implemented in an almost-truthful manner by a machine that uniformly selects $r \in \{0, 1\}^{\rho(n)}$ (and responds as the specification). However, any implementation that uses a random-tape of length ρ' , yields a function that assigns the first $\rho(n)$ arguments values that as a string have (time-bounded) Kolmogorov Complexity at most $(O(1) + \rho'(n)) + \log_2(\operatorname{poly}(n)) = \rho'(n) + O(\log n)$. Thus, for $\rho'(n) = \rho(n) - 2g(n)$, the implementation cannot even be "remotely" truthful.

Composing implementations: A simple observation that is used in our work is that one can "compose implementations". That is, if we implement a random object R1 by an oracle machine that uses oracle calls to a random object R2, which in turn has an implementation by a machine of type T, then we actually obtain an implementation of R1 by a machine of type T. To state this result, we need to extend Definition 2.5 such that it applies to oracle machines that use arbitrary specifications (rather than a random oracle). Let us denote by $(M^{(S,n)}, n)$ an implementation by the oracle machine M (and feasibility parameter n) with oracle access to the specification (S, n).

Theorem 2.17 Let $Q \in \{\text{perfect, close, pseudo}\}$. Suppose that the specification (S_1, n) can be Q-implemented by $(M^{(S_2,n)}, n)$ and that (S_2, n) has a Q-implementation by an ordinary machine (resp., by an oracle machine with a random oracle). Then, (S_1, n) has a Q-implementation by an ordinary machine (resp., by an oracle machine with a random oracle). Furthermore, if both the implementations in the hypothesis are truthful (resp., almost-truthful) then so is the implementation in the conclusion.

Proof: The idea is to simply replace (S_2, n) by its implementation, denoted (I_2, n) , and thus obtain an implementation $(M^{(I_2,n)}, n)$ of (S_1, n) , which (by combining the machines M and I_2) yields a machine of the type of machine I_2 . This machine inherits the truthfulness (resp., almost-truthfulness) of the given implementations. The analysis of the "quality" of the resulting implementation relies on the fact that the resource bounds imposed on the implementation $(M^{(S_2,n)}, n)$ induce bounds on the use of (S_2, n) by M. Combined with the hypothesis regarding the "quality" of (I_2, n) guarantees the "quality" of the resulting implementation.

For the sake of clarity, let us spell out the argument for the case of pseudo-implementations: The first hypothesis asserts that $(M^{(S_2,n)}, n)$ and (S_1, n) are computationally-indistinguishable, and the second hypothesis asserts that (I_2, n) and (S_2, n) are computationally-indistinguishable. Our goal is to prove that $(M^{(I_2,n)}, n)$ and (S_1, n) are computationally-indistinguishable, which (by the first hypothesis) reduces to proving that $(M^{(I_2,n)}, n)$ and $(M^{(S_2,n)}, n)$ are computationallyindistinguishable. Now suppose, towards the contradiction, that some a probabilistic polynomialtime machine D distinguishes $(M^{(I_2,n)}, n)$ from $(M^{(S_2,n)}, n)$. Then, combining D and M, we obtain a machine that distinguishes (I_2, n) from (S_2, n) (in contradiction to the second hypothesis). The key point is that the fact that M is probabilistic polynomialtime (because it is an implementation machine), and so the combined distinguisher is also probabilistic polynomialtime (provided that so is D). In the case of close-implementations, we rely on the fact that D makes poly(n)-many queries and each such query is served by poly(n)-many queries of M.

2.5 Objects of feasible size

In contrast to the rest of this work, we shortly discuss the complexity of generating random objects of feasible size (rather than huge random objects). In other words, we are talking about implementing a distribution on poly(n)-bit long strings, and doing so in poly(n)-time. This problem can be cast in our general formulation by considering specifications that ignore their input (i.e., have output that only depend on their random-tape). In other words, we may view objects of feasible size as constant functions, and cosider a specification of such random objects as a distribution on constant functions. Thus, without loss of generality, the implementation may also ignore its input, and consequently in this case there is no difference between an implementation by ordinary machine and an implementation by oracle machine with a random oracle.

We note that perfect implementations of such distributions were considered before (e.g., in [1, 4, 14]), and distributions for which such implementations exist are called sampleable. In the current context, where the observer sees the entire object, the distinction between perfect implementation and close-implementation seems quite technical. What seems fundamentally different is the study of pseudo-implementations.

Theorem 2.18 There exist specifications of feasible-sized objects that have no close-implementation, but do have (both truthful and non-truthful) pseudo-implementations.

Proof: Any evasive pseudorandom distribution (see [17]) yields such a specification. Recall that a distribution is called **evasive** if it is infeasible to generate an element in its support (except with negligible probability), and is called **pseudorandom** if it is computationally indistinguishable from a uniform distribution on strings of the same length. Thus, by definition, an evasive distribution has no close-implementation. On the other hand, any pseudorandom distribution can be pseudo-implemented by the uniform distribution (or any other pseudorandom distribution). Indeed, the latter implementation is not even almost-truthful (with respect to the evasive pseudorandom distribution, because even a "remotely-truthful" implementation would violate the evasiveness condition). To allow also the presentation of a truthful implementation, we modify the specification such that with exponentially-small probability it produces some sampleable pseudorandom distribution, and otherwise it produces the evasive pseudorandom distribution. The desired truthful pseudo-implementation will always produce the former distribution (i.e., the sampleable pseudorandom distribution), and still the combined distribution has no close-implementation.

The proof of Theorem 2.18 also establishes the existence of specifications (of feasible-sized objects) that have no truthful (and even no almost-truthful) implementation, regardless of the level of indistinguishability from the specification. Turning the table around, ignoring the truthfulness condition, we ask whether there exist specifications of feasible-sized objects that have no pseudo-implementations. A partial answer is provided by the following result, which relies on a non-standard assumption (see Footnote 7).

Proposition 2.19 Assuming the existence of a collision-free hash function⁷, there exists a specification of a random feasible-sized object that has no pseudo-implementation.

Proof: Given a collision-free hash function $h : \{0,1\}^{2n} \to \{0,1\}^n$, consider the uniform distribution over the set $S_n \stackrel{\text{def}}{=} \{(x,y) \in \{0,1\}^{n+n} : h(x) = h(y)\}$. Then, any implementation fails to hit the support of this distribution, which in turn is polynomial-time recognizable. Thus, the above specification (of a uniform distribution over S_n) cannot be pseudo-implemented.

Open Problem 2.20 (A stronger version of Proposition 2.19:) Provide a specification of a random feasible-sized object that has no pseudo-implementation, while relying on a standard intractability assumption.

Let us digress and consider close-implementations. For example, Bach's elegant algorithm for generating random composite numbers along with their factorization [3] can be cast as a (non-trivial) close-implementation of the said distribution.⁸ A more elementary set of examples refers to the generation of integers (out of a huge domain) according to various "nice" distributions (e.g., the binomial distribution of N trials).⁹ In fact, Knuth [25, Sec. 3.4.1] considers the generation of various such distributions, and his treatment (of integer-valued distributions) can be easily adapted to fit our formalism. This direction is further pursued in Appendix A. In general, recall that in the current context (where the observer sees the entire object), a close-implementation must be statistically close to the specification. Thus, almost-truthfulness follows "for free":

Proposition 2.21 Any close-implementation of a specification of a feasible-sized object is almosttruthful to it.

Multiple samples. Our general formulation can be used to specify an object that whenever invoked returns an independently drawn sample from the same distribution. Specifically, the specification may be by a machine that answers each "sample-query" by using a distinct portion of its random-tape (as coins used to sample from the basic distribution). Using a pseudorandom function, we may pseudo-implement multiple samples from any distribution for which one can pseudo-implement a single sample. That is:

⁹That is, for a huge $N = 2^n$, we want to generate *i* with probability $p_i \stackrel{\text{def}}{=} \binom{N}{i}/2^N$. Note $i \in \{0, 1, ..., N\}$ has feasible size, and yet the problem is not trivial (because we cannot afford to compute all p_i 's).

⁷We stress that the assumption used here (i.e., the existence of a *single* collision-free hash function) seems stronger than the standard assumption that refers to the existence of an *ensemble* of collision-free functions (cf. [9]).

⁸We mention that Bach's concrete motivation was to generate prime numbers P along with the factorization of P-1, in order to allow efficient testing of whether a given number is a primitive element modulo P. Thus, one may say that Bach's paper provides a close-implementation (by an ordinary probabilistic polynomial-time machine) of the specification that selects at random an *n*-bit long prime P and answers the query g by 1 if and only if g is a primitive element modulo P.

Proposition 2.22 Suppose that one-way functions exist, and let $D = \{D_n\}$ be a probability ensemble such that each D_n ranges over poly(n)-bit long strings. If D can be pseudo-implemented then so can the specification that answers each query by an independently selected sample of D. Furthermore, the latter implementation is by an ordinary machine and is truthful provided that the former implementation is truthful.

Proof: Consider first an implementation by an oracle machine that merely uses the random function to assign each query a random-tape to be used by the pseudo-implementation of (the single sample of the distribution) D. Since truthfulness and computational-indistinguishability are preserved by independent samples (cf. [15, Sec. 3.2.3]), we are done as far as implementations by oracle machines are concerned. Using Theorem 2.9, the proposition follows.

3 Our Main Results

We obtain several new implementations of random objects. For sake of clarity, we present the results in two categories referring to whether they yield truthful or only almost-truthful implementations.

3.1 Truthful Implementations

All implementations stated in this section are by (polynomial-time) oracle machines (which use a random oracle). Corresponding pseudo-implementations by ordinary (probabilistic polynomialtime) machines can be derived using Theorem 2.9. Namely, assuming the existence of one-way functions, each of the specifications considered below can be pseudo-implemented in a truthful manner by an ordinary probabilistic polynomial-time machine.

The basic technique underlying the following implementations is the embedding of additional structure that enables to efficiently answer the desired queries in a consistent way or to force a desired property. That is, this additional structure ensures truthfulness (with respect to the specification). The additional structure may cause the implementation to have a distribution that differs from that of the specification, but this difference is infeasible to detect (via the polynomially-many queries). In fact, the additional structure is typically randomized in order to make it undetectable, but each possible choice of coins for this randomization yields a "valid" structure (which in turn ensures truthfulness rather than only almost-truthfulness).

3.1.1 Supporting complex queries regarding boolean functions

As mentioned above, a random boolean function is trivially implemented (and in a perfect way) by an oracle machine. By this we mean that the specification and the implementation merely serve the standard evaluation queries that refer to a random function (i.e., query x is answered by the value of the function at x). Here we consider specifications that supports more powerful queries.

Example 3.1 (answering some parity queries regarding a random function): Consider a specification by a machine (and length parameter $\ell = 2n$) that, on input (i, j) where $1 \le i \le j \le 2^n$, replies with the parity of the bits in locations i through j of its random-tape. Intuitively, this machine specifies an object that, based on a random function $f : [2^n] \to \{0, 1\}$, provides the parity of the values of f on any desired interval of $[2^n]$.

Clearly, the implementation cannot afford to compute the parity of the corresponding values in its random oracle. Still, in Section 5 we present a *perfect* implementation of Example 3.1, as well as truthful close-implementations of more general types of random objects (i.e., answering any symmetric "interval" query). Specifically, we prove:

Theorem 3.2 (see Theorem 5.2)¹⁰: For every polynomial-time computable function g, there exists a truthful close-implementation of the following specification of a random object. The specification machine uses its random-tape to define a random function $f : \{0,1\}^n \to \{0,1\}$, and answers the query $(\alpha, \beta) \in \{0,1\}^{n+n}$ by $g(\sum_{\alpha < s < \beta} f(s))$.

3.1.2 Supporting complex queries regarding length-preserving functions

In Section 9 we consider specifications that, in addition to the standard evaluation queries, answer additional queries regarding a random length-preserving function. Such objects have potential applications in computational number theory, cryptography, and the analysis of algorithms (cf. [12]). Specifically, we prove:

Theorem 3.3 (see Theorem 9.2): There exists a truthful close-implementation of the following specification. The specifying machine, uniformly selects a function $f : \{0,1\}^n \to \{0,1\}^n$, and, in addition to the standard evaluation queries, answers the inverse-query $y \in \{0,1\}^n$ with the set $f^{-1}(y)$.

Alternatively, the implementation may answer with a uniformly distributed preimage of y under f (and with a special symbol in case no such preimage exists).

Theorem 3.4 (see Theorem 9.1): There exists a truthful close-implementation of the following specification. The specifying machine, uniformly selects a function $f : \{0,1\}^n \to \{0,1\}^n$, and answers the query (x,m), where $x \in \{0,1\}^n$ and $m \in [2^{\text{poly}(n)}]$, with the value $f^m(x)$ (i.e., f iterated m times on x).

This result is related to questions studied in [31, 32]; for more details, see Section 9.

3.1.3 Random graphs of various types

Random graphs have been extensively studied (cf. [6]), and in particular are known to have various properties. But does it mean that we can provide *truthful* close-implementations of uniformly distributed (huge) graphs having any of these properties?

Let us first consider a specification for a random N-vertex graph, where $N = 2^n$. Indeed, such a random graph can be specified by the machine, which viewing its random-tape ω as an N-by-N matrix, answers input $(i, j) \in [N] \times [N]$ with the value 0 if i = j, the value $\omega_{i,j}$ if i < j, and $\omega_{j,i}$ otherwise. But how about implementing a uniformly distributed graph that has various properties?

Example 3.5 (uniformly distributed connected graphs): Suppose that we want to implement a uniformly distributed connected graph (i.e., a graph uniformly selected among all connected N-vertex graph). An adequate specification may scan its random-tape, considering each N^2 -bit long portion of it as a description of a graph, and answer adjacency-queries according to the first portion that yields a connected graph. Note that the specification works in time $\Omega(N^2)$, whereas an implementation needs to work in poly(log N)-time. On the other hand, recall that a random graph is connected with overwhelmingly high probability. This suggests to implement a random connected graph by a random

¹⁰A related result was discovered before us by Naor and Reingold; see discussion in Section 5.

graph. Indeed, this yields a close-implementation, but not a truthful one (because occasionally, yet quite rarely, the implementation will yield an unconnected graph).¹¹

In Section 6 we present truthful close-implementations of Example 3.5 as well as of related specifications (i.e., of uniformly distributed graphs having various additional properties). These are all special cases of the following result:

Theorem 3.6 (see Theorem 6.2): Let Π be a monotone graph property that is satisfied by a family of strongly-constructible sparse graphs. That is, for some negligible function μ (and every N), there exists a perfect implementation of a (single) N-vertex graph with $\mu(\log N) \cdot N^2$ edges that satisfies property Π . Then, there exists a truthful close-implementation of a uniformly distributed graph that satisfies property Π .

We stress that Theorem 6.2 applies also to properties that are not satisfied (with high probability) by a random graph (e.g., having a clique of size \sqrt{N}). The proof of Theorem 6.2 relies on the following lemma, which may be of independent interest. Loosely speaking, the lemma asserts that if a monotone graph property Π is satisfied by some sparse graphs then a uniformly distributed graph having property Π is indistinguishable from a truly random graph.

Lemma 3.7 (see Lemma 6.3): Let Π be a monotone graph property that is satisfied by some N-vertex graph having $\epsilon \cdot \binom{N}{2}$ edges. Then, any machine that makes at most q adjacency queries to a graph, cannot distinguish a random N-vertex graph from a uniformly distributed N-vertex graph that satisfies Π , except than with probability $O(q\sqrt{\epsilon}) + qN^{-(1-o(1))}$.

3.1.4 Supporting complex queries regarding random graphs

Suppose that we want to implement a random N-vertex graph along with supporting, in addition to the standard adjacency queries, also some complex queries that are hard to answer by only making adjacency queries. For example suppose that on query a vertex v, we need to provide a clique of size $\log_2 N$ containing v. In Section 7 we present a truthful close-implementations of this specification:

Theorem 3.8 (see Theorem 7.2): There exists a truthful close-implementation of the following specification. The specifying machine selects uniformly an N-vertex graph and, in addition to the standard adjacency queries, answers (Log-Clique) queries of the form v by providing a random $\lceil \log_2 N \rceil$ -vertex clique that contains v (and a special symbol if no such clique exists).

Another result proved in Section 7 follows:

Theorem 3.9 (see Theorem 7.3): There exists a truthful close-implementation of the following specification. The specifying machine selects uniformly an N-vertex graph G, and in case G is Hamiltonian it uniformly selects a (directed) Hamiltonian Cycle in G, which in turn defines a cyclic permutation $\sigma : [N] \rightarrow [N]$. In addition to the standard adjacency queries, the specification answers travel queries of the form (trav, v, t) by providing $\sigma^t(v)$, and distance queries of the form (dist, v, w) by providing the smallest $t \geq 0$ such that $w = \sigma^t(v)$.

 $^{^{11}}$ Indeed, the trivial implementation (by a random graph) is almost-truthful, but here we seek a truthful implementation (because otherwise we cannot derive from it (via Theorem 2.9) even an almost-truthful pseudo-implementation by an ordinary machine).

3.1.5 Random bounded-degree graphs of various types

Random bounded-degree graphs have also received considerable attention. In Section 8 we present truthful close-implementations of random bounded-degree graphs G = ([N], E), where the machine specifying the graph answers the query $v \in [N]$ with the list of neighbors of vertex v. We stress that even implementing this specification is non-trivial if one insists on truthfully implementing *simple* random bounded-degree graphs (rather than graphs with self-loops and/or parallel edges). Furthermore, we present truthful close-implementations of random bounded-degree graphs having additional properties such as connectivity, Hamiltonicity, having logarithmic girth, etc. All these are special cases of the following result:

Theorem 3.10 (see Theorem 8.4:) Let d > 2 be fixed and Π be a graph property that satisfies the following two conditions:

- 1. The probability that Property Π is not satisfied by a uniformly chosen d-regular N-vertex graph is negligible in log N.
- 2. Property Π is satisfied by a family of strongly-constructible d-regular N-vertex graphs having girth $\omega(\log \log N)$.

Then, there exists a truthful close-implementation of a uniformly distributed d-regular N-vertex graph that satisfies property Π .

The proof relies on the following lemma, which may be of independent interest. Loosely speaking, the lemma asserts that a random isomorphic copy of a fixed *d*-regular graph of large girth is indistinguishable from a truly random *d*-regular graph.

Lemma 3.11 (see Lemma 8.1): For d > 2, let G = ([N], E) be any d-regular N-vertex graph having girth g. Let G' be obtained by randomly permuting the vertices of G (and presenting the incidence lists in some canonical order). Then, any machine M that queries the graph for the neighborhoods of q vertices of its choice, cannot distinguish G' from a random d-regular N-vertex (simple) graph, except than with probability $O(q^2/(d-1)^{(g-1)/2})$. In the case d = 2 and q < g - 1, the probability bound can be improved to $O(q^2/N)$.

3.2 Almost-Truthful Implementations

All implementations stated in this section are by *ordinary* (probabilistic polynomial-time) machines. All these results assume the existence of one-way functions.

Again, the basic technique is to embed a desirable structure, but (in contrast to Section 3.1) here the embedded structure forces the desired property only with very high probability. Consequently, the resulting implementation is only almost-truthful, which is the reason that we have to directly present implementations by ordinary machines.

A specific technique that we use is obtaining a function as a value-by-value combination of a pseudorandom function and a function of a desired combinatorial structure. The combination is done such that the combined function inherits both the pseudorandomness of the first function and the combinatorial structure of the second function (in analogy to a construction in [22]). In some cases, the combination is by a value-by-value XOR, but in others it is by a value-by-value OR with a second function that is very sparse.

3.2.1 Random codes of large distance

In continuation to the discussion in the introduction, we prove:

Theorem 3.12 (see Theorem 4.2): For $\delta = 1/6$ and $\rho = 1/9$, assuming the existence of one-way functions, there exists an almost-truthful pseudo-implementation of the following specification: The specification machine uses its random-tape to uniformly select a code $C \subset \{0,1\}^n$ having cardinality $K \stackrel{\text{def}}{=} 2^{\rho n}$ and distance at least δn , and answers the query $i \in [K]$ with the *i*-th element in C.

We comment that the above actually specifies (and implements) an encoding algorithm for the corresponding code. It would be very interesting if one can also implement a corresponding decoding algorithm; see further discussion in Section 4.

3.2.2 Random graphs of various types

Having failed to provide truthful pseudo-implementations to the following specifications, we provide almost-truthful ones.

Theorem 3.13 (see Theorem 6.5): Let $c(N) = (2 \pm o(1)) \log_2 N$ be the largest integer *i* such that the expected number of cliques of size *i* in a random N-vertex graph is larger than one. Assuming the existence of one-way functions, there exist almost-truthful pseudo-implementations of the following specifications:

- 1. A random graph of Max-Clique $c(N) \pm 1$: The specification uniformly selects an N-vertex graph having maximum clique size $c(N) \pm 1$, and answers edge-queries accordingly.
- 2. A random graph of Chromatic Number $(1 \pm o(1)) \cdot N/c(N)$: The specification uniformly selects an N-vertex graph having Chromatic Number $(1 \pm \log_2^{-1/3} N) \cdot N/c(N)$, and answers edge-queries accordingly.

Another interesting question is to provide an almost-truthful pseudo-implementation of a uniformly distributed graph having a high (global) connectivity property. Unfortunately, we do not know how to do this. Instead, we provide an almost-truthful pseudo-implementation of a random graph for which almost all pairs of vertices enjoy a high connectivity property.

Theorem 3.14 (see Theorem 6.6): For every positive polynomial p, assuming the existence of oneway functions, there exists an almost-truthful pseudo-implementation of the following specification. The specifying machine selects a graph that is uniformly distributed among all N-vertex graphs for which all but at most an $\epsilon(N) \stackrel{\text{def}}{=} 1/p(\log_2 N)$ fraction of the vertex pairs are connected by at least $(1 - \epsilon(N)) \cdot N/2$ vertex-disjoint paths. Edge-queries are answered accordingly.

Interestingly, the same implementation works for all polynomials p; that is, the implementation is independent of p, which is only used in the definition of the specification.

4 Implementing Random Codes of Large Distance

For sufficiently small $\rho, \delta > 0$, we consider codes having relative rate ρ and relative distance δ ; that is, we consider subsets $C \subset \{0,1\}^n$ such that $|C| = 2^{\rho n}$ and every two distinct codewords (i.e., $\alpha, \beta \in C$) disagree on at least δn coordinates. Such a code is called good. A random set of $K \stackrel{\text{def}}{=} 2^{\rho n}$ strings of length n is good with overwhelmingly high probability. Thus, for a random function $f:[K] \to \{0,1\}^n$, setting $C = \{f(i) : i \in [K]\}$ yields an almost-truthful close-implementation of a random code that is good, where the specification is required to answer the query i with the i-th codeword (i.e., the i-th element in the code). Recall that it is not clear what happens when we replace f by a pseudorandom function (i.e., it may be the case that the resulting code has very small distance, although most pairs of codewords are definitely far apart). To get a almost-truthful pseudo-implementation we use a different approach.

Construction 4.1 (implementing a good random code): For $k = \rho n$, we select a random k-byn matrix M, and consider the linear code generated by M (i.e., the codewords are obtained by all possible linear combinations of the rows of M). Now, using a pseudorandom function $f_s: \{0,1\}^k \rightarrow$ $\{0,1\}^n$, where $s \in \{0,1\}^n$, we consider the code $C_{M,s} = \{f_s(v) \oplus vM : v \in \{0,1\}^k\}$. That is, our implementation uses the random-tape (M,s), and provides the i-th codeword of the code $C_{M,s}$ by returning $f_s(i) \oplus iM$, where $i \in [2^k]$ is viewed as a k-dimensional row vector (or a k-bit long string).

To see that Construction 4.1 is a pseudo-implementation of a random code, consider what happens when the pseudorandom function is replaced by a truly random one (in which case we may ignore the nice properties of the random linear code generated by M).¹² Specifically, for any matrix M and any function $f:[K] \to \{0,1\}^n$, we consider the code $C_M^f = \{f(v) \oplus vM : v \in \{0,1\}^k\}$. Now, for any fixed choice of M and a truly random function $\phi:[K] \to \{0,1\}^n$, the code C_M^{ϕ} is a random code. Thus, the pseudorandomness of the function ensemble $\{f_s\}_{s \in \{0,1\}^n}$ implies that, for a uniformly chosen $s \in \{0,1\}^n$, the code $C_{M,s} = C_M^{f_s}$ is computationally indistinguishable from a random code. The reason being that ability to distinguish selected codewords of $C_M^{f_s}$ (for a random $s \in \{0,1\}^n$) from codewords of C_M^{ϕ} (for a truly random function $\phi: [K] \to \{0,1\}^n$) yields ability to distinguish the corresponding f_s from ϕ .

To see that Construction 4.1 is almost-truthful to the good code property, fix any (pseudorandom) function f and consider the code $C_M = \{f(v) \oplus vM : v \in \{0,1\}^k\}$, when M is a random k-by-n matrix. Fixing any pair of distinct strings $v, w \in \{0,1\}^k$, we show that with probability at least 2^{-3k} (over the possible choices of M), the codewords $f(v) \oplus vM$ and $f(w) \oplus wM$ are at distance at least δn , and it follows that with probability at least $1 - 2^{-k}$ the code C_M has a distance at least δn . Thus, for a random M, we consider the Hamming weight of $(f(v) \oplus vM) \oplus (f(w) \oplus wM)$, which in turn equals the Hamming weight of $r \oplus uM$, where $r = f(v) \oplus f(w)$ and $u = v \oplus w$ are fixed. The weight of $r \oplus uM$ behaves as a binomial distribution (with success probability 1/2), and thus the probability that the weight is less than δn equals $\exp(-(1 - \mathbf{H}_2(\delta)) \cdot n)$, where \mathbf{H}_2 denotes the binary entropy function. So we need $1 - \mathbf{H}_2(\delta) \cdot n > 3k$ to holds, and indeed it does hold for appropriate choices of δ and ρ (e.g, $\delta = 1/6$ and $\rho = 1/9$). Specifically, recalling that $k = \rho n$, we need $1 - \mathbf{H}_2(\delta) > 3\rho$ to hold. We get:

Theorem 4.2 For any $\delta \in (0, 1/2)$ and $\rho \in (0, 1 - \mathbf{H}_2(\delta))$, assuming the existence of one-way functions, there exists an almost-truthful pseudo-implementation by an ordinary machine of the following specification: The specification machine uses its random-tape to uniformly select a code $C \subset \{0,1\}^n$ having cardinality $K \stackrel{\text{def}}{=} 2^{\rho n}$ and distance at least δn , and answers the query $i \in [K]$ with the *i*-th element in C.

We comment that Construction 4.1 actually implements an encoding algorithm for the corresponding code, which is actually what is required in the specification. It would be very interesting if

¹²In particular, note that the resulting code is unlikely to be linear. Furthermore, any n - O(1) > k codewords are likely to be linearly independent (both when we use a random function or a pseudorandom one).

one could also implement a corresponding decoding algorithm. Note that the real challenge is to achieve "decoding with errors" (i.e., decode corrupted codewords rather than only decode uncorrupted codewords).¹³ Specifically,

Open Problem 4.3 (implementing encoding and decoding for a good random code): Provide an almost-truthful pseudo-implementation, even by an oracle machine, to the following specification. For some $\delta \in (0, 1/2)$ and $\rho \in (0, 1 - \mathbf{H}_2(\delta))$, the specification machine selects a code $C \subset \{0, 1\}^n$ as in Theorem 4.2, and answers queries of two types:

Encoding queries: For $i \in [K]$, the query (enc, i) is answered with the *i*-th element in C.

Decoding queries: For very $w \in \{0,1\}^n$ that is at distance at most $\delta n/3$ from C, the query (dec, w) is answered by the index of the (unique) codeword that is closest to w.

Indeed, we are interested in an implementation by an ordinary machine, but as stated in Section 10, it may make sense to first consider implementations by oracle machines. Furthermore, it would be nice to obtain truthful implementations, rather than almost-truthful ones. In fact, it will even be interesting to have a truthful pseudo-implementation of the specification stated in Theorem 4.2.

5 Boolean Functions and Interval-Sum Queries

In this section we show that the specification of Example 3.1 can be perfectly implemented (by an oracle machine). Recall that we seek to implement access to a random function $f: \{0,1\}^n \to \{0,1\}$ augmented with answers regarding the parity (or XOR) of the values of f on given intervals, where the intervals are with respect to the standard lex-order of n-bit string. That is, the query $q = (\alpha, \beta) \in \{0,1\}^{n+n}$, where $0^n \leq \alpha \leq \beta \leq 1^n$, is to be answered by $\bigoplus_{\alpha \leq s \leq \beta} f(s)$. The specification can answer this query in the straightforward manner, but an implementation cannot afford to do so (because a straightforward computation may take $2^n = 2^{|q|/2}$ steps). Thus, the implementation will do something completely different.¹⁴

We present an oracle machine that uses a random function $f': \bigcup_{i=0}^{n} \{0,1\}^i \to \{0,1\}$. Using f', we define $f: \{0,1\}^n \to \{0,1\}$ as follows. We consider a binary tree of depth n and associate its *i*th level vertices with strings of length i such that the vertex associated with the string s has a left (resp., right) child associated with the string s0 (resp., s1). As a mental experiment, going from the root to the leaves, we label the tree's vertices as follows:

- 1. We label the root (i.e., the level-zero vertex, which is associated with λ) by the value $f'(\lambda)$.
- 2. For i = 0, ..., n 1, and each internal vertex v at level i, we label its *left* child by the value f'(v0), and label its *right* child by the XOR of the label of v and the value f'(v0).

(Thus, the label of v equals the XOR of the values of its children.)

¹³Note that a simple modification of Construction 4.1 (e.g., replacing the *i*-th codeword, w, by the new codeword (i, w)), allows trivial decoding of uncorrupted codewords.

¹⁴The following implementation is not the simplest one possible, but we chose to present it because it generlizes to yield a proof of Theorem 5.2 (i.e., interval-sum rather than interval-sum-mod-2). A simpler implementation of Example 3.1, which does not seem to generalize to the case of interval-sum (as in Theorem 5.2), was suggested to us recently by Phil Klein, Silvio Micali and Dan Spielman. The idea is to reduce the problem of Example 3.1 to the special case where we only need to serve interval-queries for intervals starting at 0^n ; that is, we only need to serve (interval) queries of the form $(0^n, \beta)$. (Indeed, the answer to a query (α', β') , where $\alpha' \neq 0^n$, can be obtained from the answers to the queries $(0^n, \alpha'')$ and $(0^n, \beta')$, where α'' is the string preceding α' . Next observe that the query $(0^n, \beta)$ can be served by $f'(\beta)$, where $f' : \{0, 1\}^n \to \{0, 1\}$ is a random function (given as oracle).

3. The value of f at $\alpha \in \{0,1\}^n$ is defined as the label of the leaf associated with α .

By using induction on i = 0, ..., n, it can be shown that the level *i* vertices are assigned uniformly and independently distributed labels (which do depend, of course, on the level i - 1 labels). Thus, *f* is a random function. Furthermore, the label of each internal node *v* equals the XOR of the values of *f* on all leaves in the subtree rooted at *v*.

Note that the random function f' is used to directly assign (random) labels to all the left-siblings. The other labels (i.e., of right-siblings) are determined by XORing the labels of the parent and the left-sibling. Furthermore, the label of each node in the tree is determined by XORing at most n+1 values of f' (residing in appropriate left-siblings). Specifically, the label of the vertex associated with $\sigma_1 \cdots \sigma_i$ is determined by the f'-values of the strings $\lambda, 0, \sigma_1 0, ..., \sigma_1 \cdots \sigma_{i-1} 0$. Actually, the label of the vertex associated with $\alpha 1^j$, where $\alpha \in \{\lambda\} \cup \{0, 1\}^{|\alpha|-1}0$ and $j \geq 0$, is determined by the f'-values of j + 1 vertices (i.e., those associated with $\alpha, \alpha 0, \alpha 10..., \alpha 1^{j-1}0$).

$$\begin{aligned} \operatorname{label}(\alpha 1^{j}) &= \operatorname{label}(\alpha 1^{j-1}) \oplus \operatorname{label}(\alpha 1^{j-1}0) \\ &\vdots \\ &= \operatorname{label}(\alpha) \oplus \operatorname{label}(\alpha 0) \cdots \oplus \operatorname{label}(\alpha 1^{j-2}0) \oplus \operatorname{label}(\alpha 1^{j-1}0) \\ &= f'(\alpha) \oplus f'(\alpha 0) \cdots \oplus f'(\alpha 1^{j-2}0) \oplus f'(\alpha 1^{j-1}0) \end{aligned}$$

Thus, we obtain the value of f at any n-bit long string by making at most n+1 queries to f'. More generally, we can obtain the label assigned to each vertex by making at most n+1 queries to f'. It follows that we can obtain the value of $\bigoplus_{\alpha \leq s \leq \beta} f(s)$ by making $O(n^2)$ queries to f'. Specifically, the desired value is the XOR of the leaves residing in at most 2n - 1 full binary sub-trees, and so we merely need to XOR the labels assigned to the roots of these sub-trees. Actually, O(n) queries can be shown to suffice, by taking advantage on the fact that we need not retrieve the labels assigned to O(n) arbitrary vertices (but rather to vertices that correspond to roots of sub-trees with consecutive leaves). We get:

Theorem 5.1 There exists a perfect implementation (by an oracle machine) of the specification of Example 3.1.

The above procedure can be generalize to handle queries regarding any (efficiently computable) symmetric function of the values assigned by f to any given interval. In fact, it suffices to answer queries regarding the sum of these values. We thus state the following result.

Theorem 5.2 There exists a truthful close-implementation (by an oracle machine) of the following specification of a random object. The specification machine uses its random-tape to define a random function $f: \{0,1\}^n \to \{0,1\}$, and answers the query $(\alpha,\beta) \in \{0,1\}^{n+n}$ by $\sum_{\alpha < s < \beta} f(s)$.

Note that, unlike in the case of Theorem 5.1, the implementation is not perfect, which is the reason that we explicitly mention that it is truthful.

Proof: All that is needed in order to extend the "XOR construction" is to label each vertex v with the sum (rather than the sum mod 2) of the labels of all the leaves in the sub-tree rooted at v. In particular, internal nodes should be assigned random labels according to the binomial distribution, which makes the implementation more complex (even for assigning labels to the root and more so for assigning labels to left-siblings after their parents was assigned a label). Let us start with an overview:

- We label the root by a value generated according to the binomial distribution; that is, the root (of the depth-n binary tree) is assigned the value j with probability ^N_j/2^N, where N ^{def} 2ⁿ. This random assignment will be implemented using the value f'(λ), where here f' is a random function ranging over poly(n)-bit long strings rather than over a single bit (i.e., f': ∪ⁿ_{i=0} {0,1}ⁱ → {0,1}^{poly(n)}).
- 2. For i = 0, ..., n 1, and each internal vertex v at level i, we label its *left* child as follows, by using the value f'(v0). Suppose that v is assigned the value $T \leq 2^{n-i}$. We need to select a random pair of integers (l, r) such that l + r = T and $0 \leq l, r \leq 2^{n-i-1}$. Such a pair should be selected with probability that equals the probability that, conditioned on l + r = T, the pair (l, r) is selected when l and r are distributed according to the binomial distribution (of 2^{n-i-1} trials). That is, let $M = 2^{n-i}$ be the number of leaves in the tree rooted at v. Then, for l + r = T and $0 \leq l, r \leq M/2$, the pair (l, r) should be selected with probability $\binom{M/2}{l} \cdot \binom{M/2}{r} / \binom{M}{l+r}$.
- 3. As before, the value of f at $\alpha \in \{0,1\}^n$ equals the label of the leaf associated with α .

Of course, the above two types of sampling procedures have to be implemented in poly(n)-time, rather than in $poly(2^n)$ -time (and $poly(n2^{n-i})$ -time, respectively). These implementations cannot be perfect (because some of the events occur with probability $2^{-N} = 2^{-2^n}$), but it suffices to provide implementations that generates these samples with approximately the right distribution (e.g., with deviation at most 2^{-n} or so). The details concerning these implementations are provided in an Appendix A.

We stress that the sample (or label) generated for the (left sibling) vertex associated with $\alpha = \alpha' 0$ is produced based on the randomness provided by $f'(\alpha)$. However, the actual sample (or label) generated for this vertex depends also on the label assigned to its parent. (Indeed, this is different from the case of XOR.) Thus, to determine the label assigned to any vertex in the tree, we need to obtain the labels of all its ancestors (up-to the root). Specifically, let $S_1(N,\rho)$ denote the value sampled from the binomial distribution (on N trials), when the sampling algorithm uses coins ρ ; and let $S_2(T, \rho)$ denote the value assigned to the left-child, when its parent is assigned the value T, and the sampling algorithm uses coins ρ . Then, the label of the vertex associated with $\alpha = \sigma_1 \cdots \sigma_t$, denoted label(α), is obtained by computing the labels of all its ancestors as follows. First, we compute $label(\lambda) \leftarrow S_1(N, f'(\lambda))$. Next, for i = 1, ..., t, we obtain $label(\sigma_1 \cdots \sigma_i)$ by computing $label(\sigma_1 \cdots \sigma_{i-1} 0) \leftarrow S_2(label(\sigma_1 \cdots \sigma_{i-1}), f'(\sigma_1 \cdots \sigma_{i-1} 0))$, and if necessary (i.e., $\sigma_i = 1$) by computing $label(\sigma_1 \cdots \sigma_{i-1} 1) \leftarrow label(\sigma_1 \cdots \sigma_{i-1}) - label(\sigma_1 \cdots \sigma_{i-1} 0)$. That is, we first determine the label of the root (using the value of f' at λ); and next, going along the path from the root to α , we determine the label of each vertex based on the label of its parent (and the value of f' at the left-child of this parent). Thus, the computation of the label of α , only requires the value of f' on $|\alpha| + 1$ strings. As in the case of XOR, this allows to answer queries (regarding the sum of the f-values in intervals) based on the labels of O(n) internal nodes, where each of these labels depend only on the value of f' at O(n) points. (In fact, as in the case of XOR, one may show that the values of these related internal nodes depend only on the value of f' at O(n) points.)

Regarding the quality of the implementation, by the above description it is clear that the label of each internal node equals the sum of the labels of its children, and thus the implementation is truthful. To analyze its deviation from the specification, we consider the *mental experiment* in which both sampling procedures are implemented perfectly (rather than almost so), and show that in such a case the resulting implementation is perfect. Specifically, using induction on i = 0, ..., n, it can be shown that the level *i* vertices are assigned labels that are independently distributed, where each label is distributed as the binomial distribution of 2^{n-i} trials. (Indeed, the labels assigned to the vertices of level *i* do depend on the labels assigned in level i - 1.) Thus, if the deviation of the actual sampling procedures is bounded by $2^{-n} \cdot \epsilon$, then the actual implementation is at statistical distance at most ϵ from the specification.¹⁵ The latter statement is actually stronger than required for establishing the theorem.

Open problems: Theorem 5.2 provides a truthful implementation for any (feasibly-computable) symmetric function of the values assigned by a random function over any interval of $[N] \equiv \{0, 1\}^n$. Two natural extensions are suggested below.

Open Problem 5.3 (Non-symmetric queries): Provide a truthful close-implementation to the following specification. The specification machine defines a random function $f : \{0,1\}^n \to \{0,1\}$, and answers queries of the form $(\alpha,\beta) \in \{0,1\}^{n+n}$ with the value $g(f(\alpha),...,f(\beta))$, where g is some simple function. For example, consider $g(\sigma_1,...,\sigma_t)$ that returns the smallest $i \in [t]$ such that $\sigma_i \cdots \sigma_{i+\lfloor 1+\log_2 t \rfloor - 1} = 1^{1+\lfloor \log_2 t \rfloor}$ (and a special symbol if no such i exists). More generally, consider a specification machine that answers queries of the form $(k, (\alpha, \beta))$ by returning smallest $i \in [t]$ such that $\sigma_i \cdots \sigma_{i+k-1} = 1^k$, where σ_j is the j-th element in the sequence $(f(\alpha), ..., f(\beta))$.

Note that the latter specification is interesting mostly for $k \in \{\omega(\log n), ..., n + \omega(\log n)\}$. For $k \leq k_{\rm sm} = O(\log n)$ we may just make sure (in the implementation) that any consecutive interval of length $2^{k_{\rm sm}}n^2$ contains a run of $k_{\rm sm}$ ones.¹⁶ Once this is done, queries (referring to $k \leq k_{\rm sm}$) may be served (by the implementation) in a straightforward way (i.e., by scanning at most two such consecutive intervals, which in turn contain $2^{k_{\rm sm}+1}n^2 = \text{poly}(n)$ values). Similarly, for $k \geq k_{\rm lg} = n + \omega(\log n)$, we may just make sure (in the implementation) that no pair of consecutive intervals, each of length 5n, has a run of min $(k_{\rm lg}, 2n)$ ones.

Open Problem 5.4 (Beyond interval queries): Provide a truthful close-implementation to the following specification. The specification machine defines a random function $f : \{0, 1\}^n \to \{0, 1\}$, and answers queries that succinctly describe a set S, taken from a specific class of sets, with the value $\bigoplus_{\alpha \in S} f(\alpha)$. In Example 3.1 the class of sets is all intervals of $[N] \equiv \{0, 1\}^n$, represented by their pair of end-points. Another natural case is the class of sub-cubes of $\{0, 1\}^n$; that is, a set S is specified by an n-sequence over $\{0, 1, *\}$ such that the set specified by the sequence $(\sigma_1, ..., \sigma_n)$ contains the n-bit long string $\alpha_1 \cdots \alpha_n$ if and only if $\alpha_i = \sigma_i$ for every $\sigma_i \in \{0, 1\}$.

In both cases (i.e., Problems 5.3 and 5.4), even if we do not require truthfulness, the implementation may be easily distinguished from the specification if the former answers the compound queries in a non-consistent manner. At least, a potential implementation seems to be in trouble if it "lies bluntly" (e.g., answers each query by an independent random bit).

An application to streaming algorithms: Motivated by a computational problem regarding massive data streams, Feigenbaum *et. al.* [11] considered the problem of constructing a sequence of N random variables, $X_1, ..., X_N$, over $\{\pm 1\}$ such that

¹⁵We can afford to set $\epsilon = \exp(-\operatorname{poly}(n)) < 1/\operatorname{poly}(N)$, because the runing time of the actual sampling procedures is poly-logarithmic in the desired deviation.

¹⁶That is, the random function $f : [N] \to \{0, 1\}$ is modified such that, for every $j \in [N/2^{k_{sm}}n^2]$, the interval $[(j-1)2^{k_{sm}}n^2+1, ..., j2^{k_{sm}}n^2]$ contains a run of k_{sm} ones. This modification can be performed on-the-fly by scanning the relevant interval and setting to 1 a random block of k_{sm} locations if necessary. Note that, with overwhelmingly high probability, no interval is actually modified.

- 1. The sequence is "range-summable" in the sense that given $t \in [N]$ the sum $\sum_{i=1}^{t} X_i$ can be computed in poly $(\log N)$ -time.
- 2. The random variables are almost 4-wise independent (in a certain technical sense).

Using the techniques uderlying Theorem 5.2, for any $k \leq \operatorname{poly}(\log N)$ (and in particular for k = 4), we can construct a sequence that satisfies the above properties. In fact, we get sequence that is almost k-wise independent in a stronger sense than stated in [11] (i.e., we get a sequence that is statistically close to being k-wise independent).¹⁷ This is achieved by using the construction presented in the proof of Theorem 5.2, except that f' is a function selected uniformly from a family of $k \cdot (n+1)$ -wise independent functions rather than being a truly random function, where $n = \log_2 N$ (as above). Specifically, we use functions that map $\{0,1\}^{n+1} \equiv \bigcup_{i=0}^{n} \{0,1\}^{i}$ to $\{0,1\}^{\operatorname{poly}(n)}$ in a $k \cdot (n+1)$ -wise independent manner, and recall that such functions can be specified by $\operatorname{poly}(n)$ many bits and evaluated in $\operatorname{poly}(n)$ -time (since $k \leq \operatorname{poly}(n)$). In the analysis, we use the fact that the values assigned by f' to vertices in each of the (n+1) levels of the tree are k-wise independent. Thus, we can prove by induction on i = 0, ..., n, that every k vertices at level i are assigned labels according to the correct distribution (up to a small deviation). Recall that, as stated in Footnote 15, we can obtain statistical deviation that is negligible in N (in this case, with respect to a k-wise independent sequence).

A historical note: As noted above, the ideas underlying the proof of Theorem 5.2 were discovered by Moni Naor and Omer Reingold (as early as in 1999). Actually, their construction was presented within the framework of limited independence (i.e., as in the former paragraph), rather than in the framework of random functions (used throughout the rest of this section). In fact, Naor and Reingold came-up with their construction in response to a question raised by the authors of [11] (but their solution was not incorporated in [11]). The Naor–Reingold construction was used in the subsequent work of [13] (see [13, Lem. 2]). Needless to say, we became aware of these facts only after posting first versions of our work.

6 Random Graphs Satisfying Global Properties

Suppose that you want to run some simulations on huge random graphs. You actually take it for granted that the random graph is going to be Hamiltonian, because you have read Bollobas's book [6] and you are willing to discard the negligible probability that a random graph is not Hamiltonian. Suppose that you want to be able to keep succinct representations of these graphs and/or that you want to generate them using few random bits. Having also read some works on pseudorandomness (e.g., [20, 5, 33, 16]), you plan to use pseudorandom functions [16] in order to efficiently generate and store representations of these graphs. But wait a minute, are the graphs that you generate this way really Hamiltonian?

The point is that being Hamiltonian is a global property of the graph, which in turn is a huge (i.e., $\exp(n)$ -sized) object. This global property cannot be checking the adjacency of polynomially many (i.e., poly(n)-many) vertex-pairs, and so its violation cannot be translated to a contradiction of the pseudorandomness of the function. Indeed, the substitution of a random function (or a random graph) by a pseudorandom one is not guaranteed to preserve the global property. Specifically,

¹⁷This construction was actually discovered before us by Naor and Reingold (cf. [13, Lem. 2]); see further discussion at the end of this section.

it may be the case that all pseudorandom graphs are even disconnected.¹⁸ So, can we efficiently generate huge Hamiltonian graphs? As we show below, the answer to this question is positive.

In this section we consider the implementation of various types of huge random graphs. We stress that we refer to simple and labeled graphs; that is, we consider graphs without self-loops or parallel edges, and with labeled vertices (i.e., the 3-vertex graph consisting of the edge (1, 2) is different from the 3-vertex graph consisting of the edge (1, 3)). In this section, implementing a graph means answering adjacency queries; that is, the answer to the query (u, v) should indicate whether or not u and v are adjacent in the graph. Recall that the implementation ought to work in time that is poly-logarithmic in the size of the graph, and thus cannot decide "global" properties of the graph. That is, we deal with graphs having $N = 2^n$ vertices, and our procedures run in poly(n)-time.

As in Section 3, we present our results in two categories referring to whether they yield truthful or only almost-truthful implementations. In the case of truthful implementations, we show close-implementations by (polynomial-time) oracle machines (which use a random oracle), while bearing in mind that corresponding pseudo-implementations by ordinary (probabilistic polynomial-time) machines can be derived using Theorem 2.9. In contrast, in the case of almost-truthful implementations, we work directly with ordinary (probabilistic polynomial-time) machines.

6.1 Truthful implementations

Recall that a random graph (i.e., a uniformly distributed N-vertex graph) can be perfectly implemented via an oracle machine that, on input $(u, v) \in [N] \times [N]$ and access to the oracle $f : [N] \times [N] \rightarrow \{0,1\}$, returns 0 if u = v, f(u, v) if u < v, and f(v, u) otherwise. (Indeed, we merely derive a symmetric and non-reflexive version of f.)

Turning to a less trivial example, let us closely-implement a random Bipartite Graph with N vertices on each side. This can be done by viewing the random oracle as two functions, f_1 and f_2 , and answering queries as follows:

- The function f_1 is used to closely-implement a random partition of [2N] into two sets of equal size. Specifically, we use f_1 to closely-implement a permutation π over [2N], and let the first part be $S \stackrel{\text{def}}{=} \{v : \pi(v) \in [N]\}$. Let $\chi_S(v) \stackrel{\text{def}}{=} 1$ if $v \in S$ and $\chi_S(v) \stackrel{\text{def}}{=} 0$ otherwise.
- The query (u, v) is answered by 0 if $\chi_S(u) = \chi_S(v)$. Otherwise, the answer equals $f_2(u, v)$ if u < v and $f_2(v, u)$ otherwise.

The above implementation can be adapted to closely-implement a random Bipartite Graph (see details in Appendix B). Viewed in different terms, we have just discussed the implementation of random graphs satisfying certain properties.

We now turn to Example 3.5 (which specifies a uniformly distributed *connected* graph). In continuation to the discussion in Section 3, we now present a close-implementation that is *truthful*:

Construction 6.1 (Implementing a random connected graph): Use the oracle to implement a random graph, represented by the symmetric and non-reflexive random function $g : [N] \times [N] \rightarrow \{0,1\}$, as well as a permutation π over [N], which in turn is used to define a Hamiltonian path $\pi(1) \rightarrow \pi(2) \rightarrow \cdots \rightarrow \pi(N)$. Along with π , implement the inverse permutation π^{-1} , where this is

¹⁸Indeed, for each function f_s taken from some pseudorandom ensemble $\{f_s : [2^n] \times [2^n] \to \{0, 1\}\}_s$, it may hold that $f_s(v_s, u) = f_s(u, v_s) = 0$ for all $u \in [2^n]$, where v_s depends arbitrarily on f_s . For example, given a pseudorandom ensemble $\{f_s\}$ consider the ensemble $\{f_{s,v}\}$ such that $f_{s,v}(v, u) = f_{s,v}(u, v) = 0^n$ for all u's, and $f_{s,v}(x, y) = f_s(x, y)$ for all other (x, y)'s.

done by using Theorem 2.13.¹⁹ Answer the query (u, v) by 1 if and only if either g(u, v) = 1 or (u, v) is on the Hamiltonian path (i.e., $|\pi^{-1}(u) - \pi^{-1}(v)| = 1$).

Clearly, the above implementation is truthful. (Indeed, it actually implements a random Hamiltonian graph.) The implementation is statically-indistinguishable from the specification, because it is unlikely to hit an edge of the "forced Hamiltonian path" when making only poly(log N) queries. (A proof of the latter statement appears below.) A similar strategy can be used for any *monotone* graph property that satisfies the following condition:

(C) The property is satisfied by a family of strongly-constructible sparse graphs. That is, for some negligible function μ (and every N), there exists a perfect implementation of a (single) N-vertex graph with $\mu(\log N) \cdot N^2$ edges that satisfies the property.

We have:

Theorem 6.2 Let Π be a monotone graph property that satisfies Condition C. Then, there exists a truthful close-implementation (by an oracle machine) of a uniformly distributed graph that satisfies property Π .

We comment that Condition C implies that a random N-vertex graph is statistically-indistinguishable from a random N-vertex graph having property Π . This fact, which may be of independent interest, is stated and proved first.

Lemma 6.3 Let Π be a monotone graph property that is satisfied by some N-vertex graph having $\epsilon \cdot \binom{N}{2}$ edges. Then, any machine that makes at most q adjacency queries to a graph, cannot distinguish a random N-vertex graph from a uniformly distributed N-vertex graph that satisfies Π , except than with probability $O(q\sqrt{\epsilon}) + qN^{-(1-o(1))}$.

Proof: As in [19, Sec. 4], without loss of generality, we may confine ourselves to analyzing machines that inspect a random induced subgraph. That is, since both graph classes are closed under isomorphism, it suffices to consider the statistical difference between the following two distributions:

- 1. The subgraph of a uniformly distributed N-vertex graph induced by a uniformly selected set of $s \stackrel{\text{def}}{=} q + 1$ vertices.
- 2. The same vertex-induced subgraph (i.e., induced by a random set of s vertices) of a uniformly distributed N-vertex graph that satisfies property Π .

Clearly, Distribution (1) is uniform over the set of s-vertex graphs, and so we have to show that approximately the same holds for Distribution (2). Let $T \stackrel{\text{def}}{=} \binom{N}{2}$ and $M \stackrel{\text{def}}{=} \epsilon T$, and let G_0 be an N-vertex graph with M edges that satisfies property Π . Consider the set of all graphs that can be obtained from G_0 by adding $\frac{T-M}{2}$ edges. The number of these graphs is

$$\begin{pmatrix} T - M \\ \frac{T - M}{2} \end{pmatrix} = \frac{2^{T - M}}{\Theta(\sqrt{T - M})} > 2^{T - M - O(1) - \frac{1}{2} \cdot \log_2 T}$$

That is, this set contains at least a $2^{-(M+O(1)+(\log_2 T)/2)} = 2^{-\epsilon' \cdot T}$ fraction of all possible graphs, where $\epsilon' \stackrel{\text{def}}{=} \epsilon + ((\log_2 T)/2T)$. Let $X = X_1 \cdots X_T \in \{0, 1\}^T$ be a random variable that is uniformly

¹⁹That is, we use a truthful close-implementation of Example 2.4. In fact, we only need π^{-1} , and so the truthful close-implementation of Example 2.3 (as stated in Theorem 2.12) actually suffices.

distributed over the set of all graphs that satisfy property Π . Then X has entropy at least $T - \epsilon' T$ (i.e., $\mathbf{H}(X) \geq T - \epsilon' T$). It follows that $\frac{1}{T} \sum_{i=1}^{T} \mathbf{H}(X_i | X_{i-1} \cdots X_1) \geq 1 - \epsilon'$. Note that the index *i* ranges over all unordered pairs of elements of [N]. (We assume some fixed order on these pairs.) We are interested in the expected value of $\sum_{i=1}^{\binom{S}{2}} \mathbf{H}(X_{e_i(S)} | X_{e_{i-1}(S)} \cdots X_{e_1(S)})$, where $e_i(S)$ is the *i*th pair in the set $\{(u, v) : u < v \in S\}$ and S is a uniformly selected set of t vertices. Clearly

$$\mathbf{H}(X_{e_i(S)}|X_{e_{i-1}(S)}\cdots X_{e_1(S)}) \geq \mathbf{H}(X_{e_i(S)}|X_{e_i(S)-1}\cdots X_1)$$

and so

$$\mathbf{E}_{S}\left[\sum_{i=1}^{\binom{s}{2}} \mathbf{H}(X_{e_{i}(S)}|X_{e_{i-1}(S)}\cdots X_{e_{1}(S)})\right] \geq \binom{s}{2} \cdot (1-\epsilon')$$

because for a uniformly distributed $i \in [\binom{s}{2}]$ it holds that $\mathbf{E}_{S,i} \left[\mathbf{H}(X_{e_i(S)} | X_{e_i(S)-1} \cdots X_1) \right]$ equals $\mathbf{E}_j \left[\mathbf{H}(X_j | X_{j-1} \cdots X_1) \right]$, where j is uniformly distributed in [T]. Thus, for a random s-subset S, letting $Y_S = (X_{(u,v)})_{\{(u,v):u < v \in S\}}$, we have $\mathbf{E}_S[Y_S] \ge t - \epsilon''$, where $t \stackrel{\text{def}}{=} \binom{s}{2}$ and $\epsilon'' \stackrel{\text{def}}{=} t\epsilon'$. It follows (see Appendix C) that the statistical difference of Y_S from the uniform distribution over $\{0,1\}^t$ is at most $O(\sqrt{\epsilon''})$, which in turn equals $O(q\sqrt{\epsilon + T^{-(1-o(1))}})$. The lemma follows.

Proof of Theorem 6.2: Let H = ([N], E) be a graph satisfying Condition C. In particular, given $(u, v) \in [N] \times [N]$, we can decide whether or not $(u, v) \in E$ in polynomial-time. Then, using the graph H instead of the Hamiltonian path in Construction 6.1, we implement a (random) graph satisfying property Π . That is, we answer the query (u, v) by 1 if and only if either g(u, v) = 1 or (u, v) is an edge in (the "forced" copy of) H (i.e., $(\pi^{-1}(u), \pi^{-1}(v)) \in E$). Since Π is a monotone graph property, the instances of the implementation always satisfy the property Π , and thus the implementation is truthful. Furthermore, by Condition C and the fact that π is a close-implementation of a random permutation, the probability that a machine that queries the implementation for poly(log N) times hits an edge of H is negligible in log N. Thus, such a machine cannot distinguish the implementation from a random graph. Using Lemma 6.3 (with $\epsilon = \mu(\log N)$ and $q = \text{poly}(\log N)$), the theorem follows.

Examples: Indeed, monotone graph properties satisfying Condition C include Connectivity, Hamiltonicity, k-Connectivity (for every fixed k)²⁰, containing any fixed-size graph (e.g., containing a triangle or a 4-clique or a $K_{3,3}$ or a 5-cycle), having a perfect matching, having diameter at most 2, containing a clique of size at least log N, etc. All the above properties are satisfied, with overwhelmingly high probability, by a random graph. However, *Theorem 6.2 can be applied also to* (monotone) properties that are not satisfied by a random graph; a notable example is the property of containing a clique of size at least \sqrt{N} .

6.2 Almost-truthful implementations

We start by noting that if we are willing to settle for almost-truthful implementations by ordinary machines then all properties that hold (with sufficiently high probability) for random graphs can be handled easily. Specifically:

²⁰In fact, we may have $k = k(N) = \mu(\log N) \cdot N$ for any negligible function μ . The sparse graph may consist of edges between each of the N vertex and each of k(N) designated vertices.

Proposition 6.4 Let Π be any graph property that is satisfied by all but a negligible (in $\log_2 N$) fraction of the N-vertex graphs. Then, there exists an almost-truthful close-implementation (by an oracle machine) of a uniformly distributed graph that satisfies property Π .

Indeed, the implementation is by a random graph (which in turn is implemented via a random oracle). Note, however, that it is not clear what happens if we replace the random graph by a pseudorandom one (cf. Theorem 2.11). Furthermore, the proof of Theorem 2.11 can be extended to show that there exist graph properties that are satisfied by random graphs but do not have an almost-truthful implementation by an ordinary machine.²¹ In light of the above, we now focus on almost-truthful implementations by ordinary machines.

Max-clique and chromatic number. We consider the construction of pseudorandom graphs that approximately preserve the max-clique and chromatic number of random graphs.

Theorem 6.5 Let $c(N) = (2 \pm o(1)) \log_2 N$ be the largest integer *i* such that the expected number of cliques of size *i* in a random N-vertex graph is larger than one. Assuming the existence of one-way functions, there exist almost-truthful pseudo-implementations, by ordinary machines, of the following specifications:

- 1. A random graph of Max-Clique $c(N) \pm 1$: The specification uniformly selects an N-vertex graph having maximum clique size $c(N) \pm 1$, and answers edge-queries accordingly.
- 2. A random graph of Chromatic Number $(1 \pm o(1)) \cdot N/c(N)$: The specification uniformly selects an N-vertex graph having Chromatic Number $(1 \pm \log_2^{-1/3} N) \cdot N/c(N)$, and answers edge-queries accordingly.

That is, we are required to implement random-looking graphs having certain properties. Indeed, a random N-vertex graph has the above two properties with probability at least $1 - N^{-0.99}$ (cf. [6]). Thus, a random graph provides an almost-truthful close-implementation (by an oracle machine) of a uniformly selected graph having each of these properties, but it is not clear what happens when we replace the random oracle by a pseudorandom function. (In fact, one can easily construct pseudorandom functions for which the replacement yields a graph with a huge clique.) Note that Theorem 6.5 does not follow from Theorem 6.2, because the properties at hand are not monotone.²² Thus, a different approach is needed.

Proof Sketch: We start with Part 1. We define the adjacency function $g: [N] \times [N] \to \{0, 1\}$ of a graph by XORing a pseudorandom (symmetric and non-reflexive) function f with a k-wise independent function h (i.e., $g(u, v) = f(u, v) \oplus h(u, v)$), where $k \stackrel{\text{def}}{=} 5n^2$ (and $n = \log_2 N$). Recall that such k-wise independent functions can be constructed based on kn random bits. The resulting function g is both pseudorandom and k-wise independent (analogously to the construction in [22]). Now, the key observation is that the standard analysis (of the size of the max-clique in a random graph) merely refer to the expected number of cliques os fize $c(N) \pm 2$ and to its variance. Thus, this

²¹The proof of Theorem 2.11 relates to the Kolmogorov Complexity of the function (or graph). In order to obtain a graph property, we consider the minimum value of the Kolmogorov Complexity of any isomorphic copy of the said graph, and consider the set of graphs for which this quantity is greater than $N^2/4$. The latter property is satisfied by all but at most $2^{N^2/4} \cdot (N!) \ll 2^{N^2/3}$ graphs. On the other hand, the property cannot be satisfied by an instance of an implementation via an ordinary machine. Thus, any implementation (regardless of "quality") must be non-truthful (to the specification) in a strong sense.

²²For the coloring property, Condition C does not hold either.

analysis only depends on the randomness of edges within pairs of (c(N) + 2)-subsets of vertices; that is, a total of $2 \cdot \binom{c(N)+2}{2} < (c(N) + 2)^2 = (4 + o(1)) \cdot n^2$ vertex-pairs. Hence the analysis continues to hold for g (which is $5n^2$ -independent). It follows that g provides an almost-truthful pseudo-implementation of a random N-vertex graph with max-clique size $c(N) \pm 1$.

We now turn to Part 2. Let g' be the complement of a pseudorandom graph as in Part 1. We now define the adjacency function $g: [N] \times [N] \to \{0,1\}$ of a pseudorandom graph by taking the bit-wise conjunction of the pseudorandom graph g' (from above) with a function h' selected uniformly in a set H' (defined below); that is, g(u, v) = 1 iff g'(u, v) = h'(u, v) = 1. Intuitively, each function $h' \in H'$ forces a cover of [N] by N/c(N) independent sets (each of size c(N)), and so the chromatic number of g is at most N/c(N). On the other hand, each $h' \in H'$ only has independent sets of size c(N) and taking the conjunction with a random g' (which is k-wise independent for $k > {c(N)+3 \choose 2}$) is unlikely to create an independent set of size c(N)+3, and so the chromatic number of g is at least N/(c(N) + 2). Details follow.

Each function $h' \in H'$ partitions [N] to $\chi(N) = \lfloor N/c(N) \rfloor$ sets, each of size c(N), and has h'(u, v) = 1 if and only if u and v belong to different sets; that is, the complement of h' is a disjoint set of cliques each having as a vertex-set one of the sets of the partition. Thus, such h' causes each of these vertex-set to be an independent set in g. The functions in H' differ only in the partitions they use. It turns out that it suffices to use "sufficiently random" partitions. Specifically, we use $H' = \{h'_r\}_{r \in R}$, where $R = \{r \in [N] : \gcd(r, N) = 1\}$, and consider for each $r \in R$ the partition $(S_r^{(0)}, \dots, S_r^{(\chi(N)-1)}, S_r^{(\chi(N))})$, where $S_r^{(i)} = \{(c(N)i + j)r \mod N : j = 1, \dots, c(N)\}$ for $i < \chi(N)$ and $S_r^{(\chi(N))} = \{(c(N)\chi(N) + j)r \mod N : j = 1, \dots, N - c(N)\chi(N)\}$. Thus, $h'_r(u, v) = 1$ if and only if u and v do not reside in the same $S_r^{(i)}$ (i.e., $h'_r(u, v) = 0$ essentially means that $u - v \equiv jr$

(mod N) for some $j \in \{\pm(c(N) - 1)\}$). The graph G defined by g is pseudorandom because the observer is unlikely to make a query (u, v) that is affected by h'_r (because $h'_r(u, v) = 0$ yields $2(c(N) - 1) - 1 = O(\log N)$ candidates for r, which in turn is selected uniformly in the set R, where $|R| > N/O(\log N)$). The chromatic number of G is at most $\chi(N) + 1$, because its vertex-set is covered by $\chi(N) + 1$ independent sets. On the other hand, relying on the basic structure of h' and on the k-wise independence of g', we can show²³ that, with high probability, the graph G does not contain an independent set of size c(N) + 3. Thus, the chromatic number of G is at least $N/(c(N)+2) > (1-(2/c(n))\cdot\chi(N))$. Its follows that G is an almost-truthful pseudo-implementation of the desired specification.

High connectivity. Recall that in a random N-vertex graph every pair of vertices is connected by at least (1-o(1))N/2 vertex-disjoint paths. One interesting question is to provide an almost-truthful pseudo-implementation of a uniformly distributed graph having this high (global) connectivity property. Unfortunately, we do not know how to do this. A second best thing may be to provide an almost-truthful pseudo-implementation of a random graph for which almost all pairs of vertices enjoy this "high connectivity" property.

Theorem 6.6 For every positive polynomial p, assuming the existence of one-way functions, there exists an almost-truthful pseudo-implementation by an ordinary machine of the following specifi-

²³In the analysis we fix any $h' \in H'$ and show that that deleting edges as instructed by a k-wise independent function (i.e., g') is unlikely to create an independent set of size c(N) + 3. Specifically, we bound the expected number of independent set of size c(N) + 3 in the resulting graph, and thus we only rely on the independence of the selection of edges (by g') for pairs of vertices within sets of c(N) + 3 vertices. Note that the various candidate independent sets differ with respect to their intersection with the independent sets of h', and the analysis has to take this into account. The technical but elementary analysis is given in Appendix C.

cation. The specifying machine selects a graph that is uniformly distributed among all N-vertex graphs for which all but at most an $\epsilon(N) \stackrel{\text{def}}{=} 1/p(\log_2 N)$ fraction of the vertex pairs are connected by at least $(1 - \epsilon(N)) \cdot N/2$ vertex-disjoint paths. Edge-queries are answered accordingly.

Interestingly, the same implementation works for all polynomials p; that is, the implementation is independent of p, which is only needed for the definition of the specification. In fact, in contrast to all other implementations presented in this work, the implementation used in the proof of Theorem 6.6 is the straightforward one: It uses a pseudorandom function to define a graph in the obvious manner. The crux of the proof is in showing that this implementation is computationally-indistinguishable from the above specification.

Proof Sketch: We use a pseudorandom function to define a graph G = ([N], E) in the straightforward manner, and answer adjacency queries accordingly. This yields a pseudo-implementation of a truly random graph, which in turn has the strong connectivity property (with overwhelmingly high probability). Fixing a polynomial p and $\epsilon \stackrel{\text{def}}{=} \epsilon(N) \stackrel{\text{def}}{=} 1/p(\log_2 N)$, we prove that this implementation is almost-truthful to the corresponding specification. That is, we show that all but an ϵ fraction of the vertex pairs are connected via $(1 - \epsilon) \cdot N/2$ vertex-disjoint paths. We will show that if this is not the case, then we can distinguish a random graph (or function) from a pseudorandom one.

Suppose towards the contradiction that, with non-negligible probability, a pseudorandom graph violates the desired property. Fixing such a graph, G = ([N], E), our hypothesis means that at least an ϵ fraction of the vertex-pairs are connected (in G) by fewer than $(1-\epsilon) \cdot N/2$ vertex-disjoint paths. Consider such a generic pair, denoted (u, v), and define $S_0 \stackrel{\text{def}}{=} \Gamma_G(u) \cap \Gamma_G(v)$, $S_1 \stackrel{\text{def}}{=} \Gamma_G(u) \setminus \Gamma_G(v)$, and $S_2 \stackrel{\text{def}}{=} \Gamma_G(v) \setminus \Gamma_G(u)$, where $\Gamma_G(w) \stackrel{\text{def}}{=} \{x \in [N] : (w, x) \in E\}$. Note that if G were a random graph then we would expect to have $|S_0| \approx |S_1| \approx |S_2| \approx N/4$. Furthermore, we would expect to see a large (i.e., size $\approx N/4$) matching in the induced bipartite graph $B = ((S_1, S_2), E \cap (S_1 \times S_2))$; that is, the bipartite graph having S_1 on one side and S_2 on the other. So, the intuitive idea is to test that both these considition are satisfied in the pseudorandom graph. If they do then u and v are "sufficiently connected". Thus, the hypothesis that an ϵ fraction of the vertex-pairs are no "sufficiently connected" implies a distinguisher (by selecting vertex-pairs at random and testing the above two properties). The problem with the above outline is that it is not clear how to *efficiently* test that the abovementioned bipartite graph B has a sufficiently large matching.

To allow an efficient test (and thus an efficient distinguisher), we consider a more stringent condition (which would still hold in a truly random graph). We consider a fixed partition of [N]into $T \stackrel{\text{def}}{=} N/m$ parts, $(P_1, ..., P_T)$, such that $|P_i| = m = \text{poly}(n/\epsilon)$, where $n = \log_2 N$. (For example, we may use $P_i = \{(i-1)m + j : j = 1, ..., m\}$.) If G were a random graph then, with overwhelmingly high probability (i.e., at least $1 - \exp(-m^{1/O(1)}) > 1 - \exp(-n^2)$), we would have $|S_0 \cap P_i| = (m/4) \pm m^{2/3}$ for all the *i*'s. Similarly for S_1 and S_2 . Furthermore, with probability at least $1 - \exp(-n^2)$, each of the bipartite graphs B_i induced by $(P_i \cap S_1, P_i \cap S_2)$ would have a matching of size at least $(m/4) - m^{2/3}$. The key point is that we can afford to test the size of the maximium matching in such a bipartite graph, because it has 2m = poly(n) vertices.

Let us wrap-up things. If a pseudorandom graph does not have the desired property then at least ϵ fraction of its vertex-pairs are connected by less than $(1-\epsilon)N/2$ vertex-disjoint paths. Thus, samplying $O(1/\epsilon)$ vertex-pairs, we hit such a pair with constant probability. For such a vertex-pair, we consider the sets $S_{i,0} \stackrel{\text{def}}{=} P_i \cap S_0$, $S_{i,1} \stackrel{\text{def}}{=} P_i \cap S_1$ and $S_{i,2} \stackrel{\text{def}}{=} P_i \cap S_2$, for i = 1, ..., T. It must be the case that either $\epsilon/2$ fraction of the $S_{0,i}$'s are of size less than $(1 - (\epsilon/2)) \cdot (m/4)$ or that $\epsilon/2$ fraction of the bipartite subgraphs (i.e., B_i 's) induced by the pairs $(S_{1,i}, S_{2,i})$ have no matching of size $(1 - (\epsilon/2)) \cdot (m/4)$, because otherwise this vertex-pair is sufficiently connected merely by virtue of these $S_{0,i}$'s and the large matchings in the B_i 's.²⁴ We use $m > (8/\epsilon)^3$ so to guarantee that $(m/4) - m^{2/3} > (1 - (\epsilon/2))(m/4)$, which implies that (for at least an $\epsilon/2$ fraction of the *i*'s) some quantity (i.e., either $|S_{0,i}|$ or the maximum matching in B_i) is strictly larger in a random graph than in a pseudorandom graph. Now, sampling $O(1/\epsilon)$ of the *i*'s, we declare the graph to be random if all the corresponding $S_{0,i}$'s have size at least $(m/4) - m^{2/3}$ and if all the corresponding bipartite graphs B_i 's have a maximum matching of size at least $(m/4) - m^{2/3}$. Thus, we distinguish a random function from a pseudorandom function, in contradiction to the definition of the latter. The theorem follows.

Maximum Matching in most induced bipartite graphs: The proof of Theorem 6.6 can be adapted to prove the following:

Theorem 6.7 For every positive polynomial p, assuming the existence of one-way functions, there exists an almost-truthful pseudo-implementation by an ordinary machine of a uniformly selected N-vertex graph that satisfies the following property: For all but at most an $\epsilon(N) \stackrel{\text{def}}{=} 1/p(\log_2 N)$ fraction of the disjoint set-pairs $(L, R) \subseteq [N] \times [N]$ it holds that the bipartite graph induced by (L, R) has a matchning of size $(1 - \epsilon(N)) \cdot \min(|L|, |R|)$.

As in Theorem 6.6, the implementation is straightforward, and the issue is analyzing it.

Proof Sketch: Observe that almost all relevant set-pairs satisfy $|L| \approx |R| \approx N/3$, and so we focus on these pairs. It can still be shown that in a random graph, with overwhelmingly high probability, all the corresponding bipartite graphs have a sufficiently large matching. However, this will not hold if we only consider matchings that conform with the small bipartite graphs B_i 's. Still, with overwhelmingly high probability, almost all the bipartite graphs induced by pairs (L, R) as above will have a sufficiently large matching that does conform with the small bipartite graphs B_i 's. Thus, for $\epsilon = \epsilon(N)$, the distinguisher just selects $O(1/\epsilon)$ different *i*'s, and for each such *i* tests the size of the maximal matching for $O(1/\epsilon)$ random (L, R)'s. Needless to say, the distinguisher does not select such huge sets, but rather selects their projection on P_i . That is, for each such *i* (and each attempt), the distinguisher selects a random pair $(L_i, R_i) \subset P_i \times P_i$.

A different perspective: The proofs of Theorems 6.6 and 6.7 actually establish that, for the corresponding specifications, the almost-truthfulness of an implementation follows from its computational indistinguishability (w.r.t the specification).²⁵ An interesting research project is to characterize the class of specifications for which the above implication holds.

Theorem 6.8 Suppose that S is a specification for which the following two conditions hold.

²⁴That is, we get at least $((1 - (\epsilon/2)) \cdot T) \cdot ((1 - (\epsilon/2)) \cdot (m/4)) > (1 - \epsilon)(N/4)$ paths going through S_0 , and the same for paths that use the maximum matchings in the various B_i 's.

²⁵That is, these proofs establish the first condition in Theorem 6.8, whereas the second condition is established by the straightforward construction of a random graph. A key point in these examples is that, with overwhelmingly high probability, a random object in (S, n) has stronger properties that those of all objects in (S, n). This fact makes it easier to distinguish a random object in (S, n) from an object not in (S, n). For example, with overwhelmingly high probability, a random graph has larger connectivity than required in Theorem 6.6 and this connectivity is achieved via very short paths (rather than arbitrary ones). This fact enables to distinguish (S, n) from an implementation that lacks sufficiently large connectivity.

- 1. For every implementation I and every polynomial p there exists a probabilistic polynomialtime oracle machine D and a polynomial q such that if $\mathbf{Pr}[(I,n) \notin \mathrm{Supp}(S,n)] > 1/p(n)$ then $|\mathbf{Pr}[D^{(I,n)}(1^n) = 1] - \mathbf{Pr}[D^{(S,n)}(1^n) = 1]| > 1/q(n).$
- 2. S has an almost-truthful pseudo-implementation by an oracle machine that has access to a random oracle.

Then, assuming the existence of one-way function, S has an almost-truthful pseudo-implementation by an ordinary probabilistic polynomial-time machine.

Proof: Let I be the implementation guaranteed by Condition 2, and let I' be the implementation derived from I by replacing the random oracle with a pseudorandom function. Thus, I' is a pseudo-implementation of S. Using Condition 1, it follows that I' is almost-truthful to S, because otherwise we obtain an efficient oracle machine D that distinguishes I' from S.

7 Supporting Complex Queries regarding Random Graphs

In this section we provide truthful implementations of random graph while supporting complex queries, in addition to the standard adjacency queries. The graph model as in Section 6, and as in Section 6.1 we present our (truthful) implementations in terms of oracle machines. Let us start with a simple example.

Proposition 7.1 There exists a truthful close-implementation by an oracle machine of the following specification. The specifying machine selects uniformly an N-vertex graph and answers distance queries regarding any pair of vertices. Furthermore, there exists a truthful close-implementation of the related specification that returns a uniformly distributed path of shortest length.

Proof: Consider the property of having diameter at most 2. This property satisfies Condition C (e.g., by an N-vertex star). Thus, using Theorem 6.2, we obtain a close-implementation of a random graph, while our implementation always produces a graph having diamater at most 2 (or rather exactly 2). Now, we answer the query (u, v) by 1 if the edge (u, v) is in the graph, and by 2 otherwise. For the furthermore-part, we add \sqrt{N} such stars, and serve queries regarding paths of length 2 by using the center of one of these stars (which is selected by applying an independent random function to the query pair).

This example is not very impressive because the user could have served the distance-queries in the same way (by only using adjacency queries to the standard implementation of a random graph). (A random shortest path could have also been found by using the standard implementation.) The only advantage of Proposition 7.1 is that it provides a truthful implementation of the distancequeries (rather than merely an almost-truthful one obtained via the trivial implementation). A more impressive example follows. Recall that a random N-vertex graph is likely to have many $(\log_2 N)$ -vertex cliques that include each of the vertices of the graph, whereas it seems hard to find such cliques (where in *hard* we mean unlikely to achieve in time poly(log N), and not merely in time poly(N)). Below we provide an implementation of a service that answers queries of the form $v \in [N]$ with a log-sized clique containing the vertex v.

Theorem 7.2 There exists a truthful close-implementation of the following specification. The specifying machine selects uniformly an N-vertex graph and, in addition to the standard adjacency queries, answers (Log-Clique) queries of the form v by providing a random $\lceil \log_2 N \rceil$ -vertex clique that contains v (and a special symbol if no such clique exists).

Proof Sketch: Let $\ell = \lceil \log_2 N \rceil - 1$ and consider a simple partition of [N] to $T = \lceil N/\ell \rceil$ subsets, $S_1, ..., S_T$, such that $|S_i| = \ell$ for i = 1, ..., T - 1 (e.g., $S_i = \{(i - 1)\ell + j : j = 1, ..., \ell\}$). Use the oracle to implement a random graph, G' = ([N], E'), as well as a random $onto^{26}$ function $f : [N] \to [T]$ and a random invertible permutation $\pi : [N] \to [N]$ (as in Theorem 2.13). The graph we implement will consist of the union of G' with N cliques, where the *i*-th clique resides on the vertex set $\{i\} \cup \{\pi(j) : j \in S_{f(i)}\}$. The Log-Clique queries are served in the obvious manner; that is, query v is answered with $\{v\} \cup \{\pi(u) : u \in S_{f(v)}\}$. (For simplicity, we ignore the unlikely case that $v \in \{\pi(u) : u \in S_{f(v)}\}$; this can be redeemed by letting $\ell = \lceil \log_2 N \rceil$ and answering with a random ℓ -subset of $\{v\} \cup \{\pi(u) : u \in S_{f(v)}\}$ that contains v.) Implementing the adjacency queries is slightly more tricky. The query (u, v) is answered by 1 if and only if either $(u, v) \in E$ or u and v reside in one of the N's cliques we added. The latter case may happen if and only if one of the following subcases holds:

- 1. Either $u \in \{\pi(w) : w \in S_{f(v)}\}$ or $v \in \{\pi(w) : w \in S_{f(u)}\}$; that is, either $\pi^{-1}(u) \in S_{f(v)}$ or $\pi^{-1}(v) \in S_{f(u)}$. Each of these conditions is easy to check by invoking f and π^{-1} .
- 2. There exists an x such that $u, v \in \{\pi(w) : w \in S_{f(x)}\}$, which means that $\pi^{-1}(u), \pi^{-1}(v) \in S_{f(x)}$. Equivalently, recalling that f is onto, we may check whether there exists a y such that $\pi^{-1}(u), \pi^{-1}(v) \in S_y$, which in turn is easy to determine using the simple structure of the sets S_y 's (i.e., we merely tests whether or not $\lceil \pi^{-1}(u)/\ell \rceil = \lceil \pi^{-1}(v)/\ell \rceil$).

Thus, our implementation is truthful to the specification. To see that it is a close-implementation of the specification, observe that it is unlikely that two different Log-Clique queries are "served" by the same clique (becuase this means forming a collision under f). Conditioned on this rare event not occurring, the Log-Clique queries are served by disjoint random cliques, which is what would essentially happen in a random graph (provided that poly(log N) queries are made). Finally, it is unlikely that the answers to the adjacency queries that are not determined by prior Log-Clique queries be affected by the sparse sub-graph (of N small cliques) that we inserted under a random permutation. The theorem follows.

Another example: We consider the implementation of a random graph along with answering queries regarding a random Hamiltonian cycle in it, where such cycle exists with overwhelmingly high probability. Specifically, we consider queries of the form *what is the distance between two vertices on the cycle*.

Theorem 7.3 There exists a truthful close-implementation of the following specification. The specifying machine selects uniformly an N-vertex graph G, and in case G is Hamiltonian it uniformly selects a (directed) Hamiltonian Cycle in G, which in turn defines a cyclic permutation $\sigma : [N] \rightarrow [N]$. In addition to the standard adjacency queries, the specification answers travel queries of the form (trav, v, t) by providing $\sigma^t(v)$, and distance queries of the form (dist, v, w) by providing the smallest $t \geq 0$ such that $w = \sigma^t(v)$.

We stress that the implementation must answer each possible query in time polynomial in the vertex name (which may be logarithmic in the distance t).

²⁶Such a function may be obtained by combining the identity function over [T] with a random function f': $\{T + 1, ..., N\} \rightarrow [T]$, and randomly permuting the domain of the resulting function.

Proof Sketch: It will be convenient to use the vertex set $V = \{0, 1, ..., N - 1\}$ (instead of [N]). We use the random oracle to implement a random graph G' = (V, E') as well as a random permutation $\pi : V \to V$ along with its inverse. We define a graph G = (V, E) by $E \stackrel{\text{def}}{=} E' \cup C$, where $C = \{(\pi(i), \pi(i+1 \mod N)) : i \in V\}$, and use C to answer the special (Hamiltonian) queries. That is, we answer the query (trav, v, t) by $\pi(\pi^{-1}(v) + t \mod N)$, and the query (dist, v, w) by $\pi^{-1}(w) - \pi^{-1}(v) \mod N$. The standard adacency query (u, v) is answered by 1 if and only if either $(u, v) \in E \text{ or } \pi^{-1}(u) \equiv \pi^{-1}(v) \pm 1 \pmod{N}$. (Indeed, the above construction is reminiscent of the "fast-forward" construction of [31] (stated in Theorem 2.14).)

To see that the above truthful implementation is statistically-indistinguishable from the specification, we use the following three observations:

- 1. If a (labeled) graph appears in the specification (resp., in the implementation) then all is (labeled) isomorphic copies appear in it. Consequently, for any Hamiltonian Cycle, the set of Hamiltonian graphs in which this cycle has been selected in the specification (resp., in the implementation) is isomorphic to the set of Hamiltonian graphs in which any other Hamiltonian cycle has been selected. Thus, we may consider the conditional distribution induced on the specification (resp., on the implementation) by fixing any such Hamiltonian Cycle.
- 2. Conditioned on any fixing Hamiltonian Cycle being selected in the implementation, the rest of the graph selected by the implementation is truly random.
- 3. Conditioned on any fixing Hamiltonian Cycle being selected in the specification, the rest of the graph selected by the specification is indistinguishable from a random graph. The proof of this assertion is similar to the proof of Lemma 6.3. The key point is proving that, conditioned on a specific Hamiltonian Cycle being selected, the (rest of the) graph selected by the specification has sufficiently high entropy. Note that here we refer to the entropy of the remaining $\binom{N}{2} N$ edges, and that the vertex pairs are not all identical but rather fall into categories depending on their distance as measured on the selected Hamiltonian Cycle. We need to show that a random vertex-pair in *each* of these categories has a sufficiently high (conditional) entropy. Thus, this observation requires a careful proof to be presented next.

Indeed, the above discussion suggests that we may give the entire Hamiltonian cycle to the machine that inspects the rest of the graph (in an attempt to distinguish the implementation from the specification). Thus, we assume, without loss of generality, that this machine makes no adjacency queries regarding edges that participate in the cycle. The first observation says that we may consider any fixed cycle, and the second observation says that a machine that inspects the rest of the graph sees truly random edges. The third observation, proved below, asserts that making a few queries to the rest of the conditional space of the specification, yields answers that also look random.

We consider the conditional distribution of the rest of the graph selected by the specification, given that a specific Hamiltonian Cycle was selected. (Indeed, we ignore the negligible (in N) probability that the graph selected by the specification is not Hamiltonian.) Using Bayes' Law, the conditional probability that a specific graph is selected is inversely proportional to the number of Hamiltonian Cycles in that graph. Using known results on the concentration of the latter number in random graphs (see, e.g., [24, Thm. 4]), we infer that in all but an N^{-2} fraction of the N-vertex graphs the number of Hamiltonian Cycles is at least an $\exp(-2(\ln N)^{1/2}) > N^{-1}$ fraction of its expected number. Thus, the conditional entropy of the selected graph (conditioned on the selected cycle) is $\binom{N}{2} - N - o(N)$. Details follow.

For $T = \binom{N}{2}$, let $X = X_1 \cdots X_T$ denote the graph selected by the specification, and Y(G) denote the Hamiltonian Cycle selected (by the specification) given that the graph G was selected. Let $\#_{\mathrm{HC}}(G)$ denote the number of Hamiltonian Cycles in the graph G, where cyclic shifts and traspositions of cycles are counted as if they were different cycles (and so the number of Hamiltonian Cycles in an N-clique is N!). Thus, $\mathbf{E}(\#_{\mathrm{HC}}(X)) = 2^{-N} \cdot (N!)$. An N-vertex graph G is called good if $\#_{\mathrm{HC}}(G) > 2^{-N} \cdot (N-1!)$, and \mathcal{G} denotes the set of good N-vertex graphs. For a Hamiltonian Cycle C, we denote by $\mathcal{G}(C)$ the set of graphs in \mathcal{G} that contain the cycle C. Then, it holds that

$$\begin{split} \mathbf{H}(X|Y(X) &= C) &\geq \sum_{G \in \mathcal{G}(C)} \mathbf{Pr}[X = G|Y(X) = C] \cdot \log_2(1/\mathbf{Pr}[X = G|Y(X) = C])) \\ &\geq (1 - N^{-2}) \cdot \min_{G \in \mathcal{G}(C)} \left\{ -\log_2(\mathbf{Pr}[X = G|Y(X) = C])) \right\} \\ &= (1 - N^{-2}) \cdot \min_{G \in \mathcal{G}(C)} \left\{ \begin{array}{l} \log_2(\mathbf{Pr}[Y(X) = C]) \\ -\log_2(\mathbf{Pr}[Y(X) = C|X = G]) \\ -\log_2(\mathbf{Pr}[X = G]) \end{array} \right\} \\ &= (1 - N^{-2}) \cdot \min_{G \in \mathcal{G}(C)} \left\{ \log_2(1/N!) + \log_2(\#_{\mathrm{HC}}(G)) + \binom{N}{2} \right\} \end{split}$$

Using the fact that G is good (i.e., $G \in \mathcal{G}(C)$), it follows that $\log_2(\#_{\mathrm{HC}}(G)) > \log_2(2^{-N} \cdot (N-1!))$, which in turn equals $\log_2(N!) - N - \log_2 N$. We thus get,

$$\mathbf{H}(X|Y(X) = C) > (1 - N^{-2}) \cdot \left(\binom{N}{2} - N - \log_2 N \right)$$
(2)

Recall that the condition Y(X) = C determines N vertex-pairs in X, and so the entropy of the remaining $T' = \binom{N}{2} - N$ pairs is at least $T' - \log_2 N$. Partitioning these (undetermined) pairs according to their distances in C, we conclude that the entropy of the N/2 pairs in each such distance-class is at least $(N/2) - \log_2 N$. (Indeed, the distance class of undetermined pairs do not contain distance 1 (or N - 1), which correspond to the forced cycle-edges.) We stress that our analysis holds even if the machine inspecting the graph, is given the Hamiltonian cycle for free. This machine may select the induced subgraph that it wants to inspect, but this selection is determined upto a shifting of all vertices (i.e., a rotation of the cycle). This randomization suffices for concluding that the expected entropy of the inspected subgraph (which may not include cycle edges) is at least $(1 - ((2 \log_2 N)/N)) \cdot \binom{t}{2})$, where t is the number of vertices in the subgraph. As in the proof of Lemma 6.3, this implies that the inspected subgraph is at distance at most $O(\sqrt{((\log_2 N)/N) \cdot \binom{t}{2}}) < t \cdot N^{-(1-o(1))/2}$ from a random t-vertex graph. The theorem follows.

8 Random Bounded-Degree Graphs and Global Properties

In this section we consider huge bounded-degree simple graphs, where the vertices are labelled (and there are no self-loops or parallel edges). We consider specifications of various distributions over such graphs, where in all cases the specifying machine responds to neighborhood queries (i.e., the queries correspond to vertices and the answer to query v is the list of vertices that are adjacent to vertex v).

The first issue that arises is whether we can implement a random bounded-degree graph or alternatively a random regular graph. Things would have been quite simple if we were allowing also non-simple graphs (i.e., having self-loops and parallel edges). For example, a random *d*-regular *N*-vertex non-simple graph can be implemented by pairing at random the dN possible "ports" of the *N* vertices. We can avoid self-loops (but not parallel edges) by generating the graph as a union of *d* perfect matchings of the elements in [N]. In both cases, we would get a close-implementation of a random *d*-regular *N*-vertex (simple) graph, but parallel edges will still appear with constant probability (and thus this implementation is not truthful w.r.t simple graphs). In order to obtain a random simple *d*-regular *N*-vertex graph, we need to take an alternative route. The key observation underlying this alternative is captured by the following lemma:

Lemma 8.1 For d > 2, let G = ([N], E) be any d-regular N-vertex graph having girth g. Let G' be obtained by randomly permuting the vertices of G (and presenting the incidence lists in some canonical order). Then, any machine M that queries the graph for the neighborhoods of q vertices of its choice, cannot distinguish G' from a random d-regular N-vertex (simple) graph, except than with probability $O(q^2/(d-1)^{(g-1)/2})$. In the case d = 2 and q < g-1, the probability bound can be improved to $O(q^2/N)$.

Recall that the girth of a graph G is the length of the shortest simple cycle in G, and that $(d - 1)^{(g-2)/2} < N$ always holds (for a d-regular N-vertex graph of girth g).²⁷ Note that Lemma 8.1 is quite tight: For example, in the case d = 2, for $g \ll \sqrt{N}$, the N-vertex graph G may consist of a collection of g-cycles, and taking a walk of length g in G' (by making g - 1 queries) will always detect a cycle G', which allows to distinguish G' from a random 2-regular N-vertex (in which the expected length of a cycle going through any vertex is $\Omega(N)$). In the case d > 3, the graph G may consist of connected components, each of size $(d - 1)^g \ll N$, and taking a random walk of length $(d - 1)^{g/2}$ in G' is likely to visit some vertex twice, which allows to distinguish G' from a random d-regular N-vertex (in which this event may occur only after \sqrt{N} steps). Below, we will use Lemma 8.1 with the following setting of parameters.

Corollary 8.2 For fixed d > 2 and $g(N) = \omega(\log \log N)$, let G = ([N], E) be any d-regular N-vertex graph having girth g(N). Let G' be obtained from G as in Lemma 8.1. Then, any machine M that queries the graph for the neighborhoods of poly $(\log N)$ vertices of its choice, cannot distinguish G' from a random d-regular N-vertex (simple) graph, except than with negligible in $\log N$ probability. The claim holds also in the case that d = 2 and $g(N) = (\log N)^{\omega(1)}$.

For d > 2 the girth can be at most logarithmic, and explicit constructions with logarithmic girth are known for all $d \ge 3$ and a dense set of N's (which is typically related to the set of prime numbers; see, e.g., [30, 23, 28]). For d = 2, we may just take the N-cycle or any N-vertex graph consisting of a collection of sufficiently large cycles.

Proof Sketch for Lemma 8.1: We bound the distinguishing gap of an oracle machine (which queries either a random *d*-regular *N*-vertex graph or the random graph G') as a function of the number of queries it makes. Recall that G' is a random isomorphic copy of G, whereas a random *d*-regular *N*-vertex graph may be viewed as a random isomorphic copy of another random *d*-regular *N*-vertex graph. Thus, intuitively, the specific labels of queried vertices and the specific labels of the corresponding answers are totally irrelevant: the only thing that matters is whether or not

²⁷The girth upper-bound (i.e., $g \leq 2 + 2\log_{d-1} N$) follows by considering the (vertex disjoint) paths of length (g-2)/2 starting at any fixed vertex. The existence of *d*-regular *N*-vertex graphs of girth $\log_{d-1} N$ was shown (non-constructively) in [10].

two labels are equal.²⁸ Equality (between labels) can occur in two cases. The uninteresting case is when the machine queries a vertex u that is a neighbor of a previously-queried vertex v and the answer contains (of course) the label of vertex v. (This is uninteresting because the machine, having queried v before, already knows that v is a neighbor of u.) The interesting case is that the machine queries a vertex and the answer contains the label of a vertex v that was not queried before but has already appeared in the answer to a different query. An important observation is that, as long as no interesting event occurs, the machine cannot distinguish the two distributions (becuase in both cases it knows the same subgraph, which is a forest). Thus, the analysis amounts to bounding the probability that an interesting event occurs, when we make q queries.

Let us consider first what happens when we query a random d-regular N-vertex (simple) graph. We may think of an imaginary process that constructs the graph on-the-fly such that the neighbors of vertex v are selected only in response to the query v (cf. [18, Thm. 7.1]). This selection is done at random according to the conditional distribution that is consistent with the partial graph determined so far. It is easy to see that the probability that an interesting event occurs in the *i*-th query is at most (i - 1)d/(dN - (i - 1)d), and so the probability for such an event occurring in qqueries is at most q^2/N .

The more challenging part is to analyse what happens when we query the graph G. (Recall that we have already reduced the analysis to a model in which we ignore the specific labels, but rather only compare them, and analogously we cannot query a specific new vertex but rather only query either a random new vertex or a vertex that has appeared in some answer.)²⁹ To illustrate the issues at hand, consider first the case that d = 2 (where G consists of a set of cycles, each of length at least g). In this case, we have the option of either to proceed along a path that is part of a cycle (i.e., query for the neighbors of the an end-point of a currently known path) or to query for a random new vertex. Assuming that we make less than g - 1 queries, we can never cause an interesting event by going along a path (because an interesting event may occur in this case only if we go around the entire cycle, which requires at least g - 1 queries). The only other possibility to encounter an interesting event is by having two paths (possiblly each of length 1) collide. But the probability for such an event is bounded by q^2/N , where q is the number of queries that we make.³⁰

We now turn to the more interesting case of d > 2. As in case d = 2, taking a walk of length g - 2 from any vertex will not yield anything useful. However, in this case, we may afford to take longer walks (because q may be much larger than g). Still, we will prove that, in this case, with probability at least $1 - q^2 \cdot (d - 1)^{-(g-3)/2}$, the uncovered subgraph is a forest. The proof relies both on the the girth lower-bound of G and on a sufficiently-good rapid-mixing property (which follows from the girth lower-bound). We bound the probability that a cycle is closed in the current forest by the probability that two vertices in the forest are connected by a non-tree edge, where the probability is taken over the possible random vertices returned in response to a new-vertex request and over the random order in which neighbors of a query-vertex are provided.

²⁸Essentially, the machine cannot determine which vertex it queries; all that it actually decides is whether to query a specific vertex that has appeared in previous answers or to query a new vertex (which may be viewed as randomly selected). (Formally, a specific new label indicated by the querying machine is mapped by the random permutation to a new random vertex.) Similarly, the labels of the vertices given as answer do not matter, all that matters is whether or not these vertices have appeared in the answers to previous queries (or as previous queries). (Again, formally, the new vertices supplied in the answer are assigned, by the random permutation, new random labels.)

²⁹Thus, we may consider querying G itself (rather than querying G').

³⁰Using a union bound over all query pairs, we bound the probability that the *i*th query collides with the *j*-th query. Each of these two queries is obtained by a path of fixed length starting from a uniformly and distributed vertex (which was new at the time). Thus, these two queries are almost uniformly and independently distributed (in [N]), and the probability that they are neighbors is at most 1/(N-q).

Indeed, a key observation is that when we query a vertex that has appeared in some answer, we may think that this vertex is selected at random among the unqueried vertices appearing in that answer.³¹ Taking a union bound on all possible $\binom{q}{2}$ vertex pairs (i.e., those in the forest), we bound the probability that either two ends of a discovered path (in one tree) or two vertices in different current trees are connected by an edge. (In both cases, these vertices are actually leaves.)

We consider each of these two cases seperately: In the latter case (i.e., leaves in different trees), the two vertices (which are not connected in the currently uncovered subgraph) are uniformly distributed in G, and thus the probability that they are connected is essentially d/N. The situation here is essentially as analyzed in the case d = 2: we have two paths, each initiated at a random (new at the time) vertex, leading to the leaves in question, and thus the latter are almost uniformly and independently distributed.

Turning to the former case (i.e., endpoints of a path in a tree), we use the girth hypothesis to infer that this path must have length at least g-1 (or else its endpoint are definitely not connected). However, the machine that discovered this path actually took a random walk (possiblly to two directions) starting from one vertex, because we may assume that this is the first time in which two vertices in the current forest are connected by a current non-tree edge. We also use the hypothesis that our exploration of the path (i.e., queries regarding vertices that appeared in previous answers) is actually random (i.e., we effectively extend the current end-point of the path by a uniformly selected neighbor of that end-point). Now, the end-point of such a path cannot hit any specific vertex with probability greater than $\nu \stackrel{\text{def}}{=} (d-1)^{-(g-1)/2}$, because after (g-1)/2 steps the end-point must be uniformly distributed over the $(d-1)^{(g-1)/2}$ leaves of the tree rooted at the start vertex (and the max-norm of a distribution cannot increase by additional random steps). Fixing the closest (to the start vertex) end-point, it follows that the probability that the other end-point hits the neighbor-set of the first end-point is at most $d \cdot \nu = O((d-1)^{-(g-1)/2})$. To summarize, the probability that an interesting event occurs, while making q queries, is at most $O(q^2 \cdot (d-1)^{-(g-1)/2})$. The lemma follows.

Implementing random bounded-degree simple graphs: We now turn back to the initial problem of implementing random bounded-degree (resp., regular) simple graphs.

Proposition 8.3 For every constant d > 2, there exist truthful close-implementations of the following two specifications:

- 1. A random graph of maximum degree d: For size parameter N, the specification selects uniformly a graph G among the set of N-vertex simple graphs having maximum degree d. On query $v \in [N]$, the machine answers with the list of neighbors of vertex v in G.
- 2. A random d-regular graph: For size parameter N, the specification selects uniformly a graph G among the set of N-vertex d-regular simple graphs, and answers queries as in Part 1.

Proof: We start with Part 2. This part should follow by Corollary 8.2, provided that we can implement a random isomorphic copy of a *d*-regular *N*-vertex graph of sufficiently large girth.

³¹That is, the correspondance between the new place-holders in the answer and the new real neighbors of the queried vertex is random. Formally, we may define the interaction with the graph such that at each point only the internal nodes of the currently revealed forest are assigned a serial number. Possible queries may be either for a new random vertex (assigned the next serial number and typically initiating a new tree in the forest) or for a random leaf of a specific internal vertex (which typically extends the corresponding tree and turns one of these leaves to an internal vertex with d-1 new leaves).

This requires an explicit construction of the latter graph as well as an implementation of a random permutation and its inverse (as provided by Theorem 2.13). Specifically, let G_N be the fixed graph, and π the random relabelling of its vertices. The we answer query v, by first determining the preimage of v in G_N (i.e., $\pi^{-1}(v)$), next find its neighbors (using the explicitness of the construction of G_N), and finally return their images under π . Indeed, this process depends on the ability to provide explicit constructions of adequate d-regular N-vertex graphs (i.e., G_N 's). This is trivial in the case d = 2 (e.g., by the N-cycle). For other values of $d \geq 3$, adequate constructions can be obtained from [30, 23, 28, 26] (possibly by dropping several (easily identified) perfect matchings from the graph). These construction apply for a dense set of N's (which are typically of the form $p(p-1)^2$ for any prime p), but we can obtain other sizes by combining many such graphs (note that we are not even required to give a connected graph, let alone a good expander).

We now turn to Part 1. We first note that most graphs of maximum degree d have $(1-o(1)) \cdot dN/2$ edges. Furthermore, for $T = \Theta(\sqrt{dN})$ and $D = O(\sqrt{dN})$, all but a negligible (in N) fraction of the graphs have $(dN/2) - T \pm D$ edges. In this range, random N-vertex graphs with a given number of edges and degree bound d, can be closely-implemented by selecting a random d-regular N-vertex graph and omitting the adequate number of edges. Thus, all that is needed is to select M at random with probability proportional to the number of N-vertex graphs with M edges and degree bound d. This can be done by using known expressions for these numbers, and techniques such as in Appendix A.

A general result: The proof of Proposition 8.3 actually yields a truthful close-implementation of several other specifications. Consider, for example, the generation of random connected *d*-regular graphs, for $d \ge 3$. Since the explicit constructions of *d*-regular graphs are connected (and their modifications can easily made connected), applying Corollary 8.2 will do. (Indeed, we also use the fact that, with overwhelmingly high probability, a random *d*-regular graph is connected.) More generally, we have:

Theorem 8.4 Let $d \ge 2$ be fixed and Π be a graph property that satisfies the following two conditions:

- 1. The probability that Property Π is not satisfied by a uniformly chosen d-regular N-vertex graph is negligible in $\log N$.
- 2. Property Π is satisfied by a family of strongly-constructable d-regular N-vertex graphs having girth $\omega(\log \log N)$ if d > 2 and girth $(\log N)^{\omega(1)}$ if d = 2.

Then, there exists a truthful close-implementation (by an oracle machine) of a uniformly distributed d-regular N-vertex graph that satisfies property Π .

We note that Condition 1 may be relaxed. It suffices to require that a random d-regular graph and a random d-regular graph having Property II are statistically-indistinguishable (by a machine that makes poly-logarithmically many queries). In particular, a random 2-regular graph and a uniformly distributed *connected* 2-regular graph are statistically-indistinguishable, and thus we can provide a truthful close-implementation of the latter specification. We mention that Theorem 8.4 yields truthful close-implementations to random d-regular graphs that are required to be Hamiltonian, Bipartite, have logarithmic girth, etc.

9 Supporting Complex Queries regarding Length-Preserving Functions

In this section we consider specifications that, in addition to the standard evaluation queries, answer various queries regarding a random function $f : \{0,1\}^n \to \{0,1\}^n$. The first type of queries we handle are interated-evaluation queries, where the number of iterations may be super-polynomial in the length of the input (and thus cannot be implemented in a straightforward manner).

Theorem 9.1 There exists a truthful close-implementation of the following specification. The specifying machine, uniformly selects a function $f : \{0,1\}^n \to \{0,1\}^n$, and answers queries of the form (x,m), where $x \in \{0,1\}^n$ and $m \in [2^{poly(n)}]$, with the value $f^m(x)$ (i.e., f iterated m times on x).

Proof: Consider first an implementation by a random N-cycle, where $N = 2^n$. That is, using a random 1-1 mapping $\pi : \{0, ..., N-1\} \rightarrow \{0, 1\}^n$, define $f(x) = \pi(\pi^{-1}(x) + 1 \mod N)$, and answer the query (x, m) by $\pi(\pi^{-1}(x) + m \mod N)$. (Indeed, the above construction is reminiscent of the "fast-forward" construction of [31] (stated in Theorem 2.14).) The only thing that goes wrong is that we know the cycle length of f and thus can distinguish it from a random function by any query of the form (\cdot, N) . Thus, we modify the construction so to obtain a function f with unknown cycle lengths. A simple way of doing this is to use two cycles, while randomly selecting the length of the first cycle. That is, select M uniformly in [N], and let

$$f(x) \stackrel{\text{def}}{=} \begin{cases} \pi(\pi^{-1}(x) + 1 \mod M) & \text{if } \pi^{-1}(x) \in \{0, ..., M - 1\} \\ \pi(\pi^{-1}(x) + 1) & \text{if } \pi^{-1}(x) \in \{M, ..., N - 2\} \\ \pi(M) & \text{otherwise (i.e., } \pi^{-1}(x) = N - 1) \end{cases}$$

We could have tried to select f such that its cycle structure is distributed as in case of a random function, but we did not bother to do so. Nevertheless, we prove that any machine that makes q queries cannot distinguish f from a random function with probability better than $poly(n) \cdot q^2/2^{\Omega(n)}$. Actually, in order to facilitate the analysis, we select M uniformly in $\{(N/3), ..., (2N/3)\}$.

We turn to prove that the above truthful implementation is statistically-indistinguishable from the specification. As in the proof of Lemma 8.1, we may disregard the actual values of queries and answers in the querying process, and merely refer to whether these values are equal or not. We also assume, without loss of generality, that the querying machine makes no redundent queries (e.g., if it knows that $y = f^k(x)$ and $z = f^{\ell}(y)$ then it refrains from making the query $(x, k + \ell)$, which would have been answered by $z = f^{k+\ell}(x)$). That is, at any point in time, the querying machine knows of a few chains, each having the form $(x, f^{k_1}(x), f^{k_2}(x), ..., f^{k_t}(x))$, for some known $x \in \{0, 1\}^n$ and $k_1 < k_2 < \cdots < k_t$. Typically, the elements in each chain are distinct, and no element appears in two chains. In fact, as long as this typical case holds, there is no difference between querying the specification versus querying the implementation. Thus, we have to upper bound the probability that an untypical event occurs (i.e., a query is answered by an element that already appears on one of the chains, although the query was not redundent).

Let us first consider the case that f is constructed as in the implementation. For the *i*-th non-redundent query, denoted (x, k), we consider three cases:

Case 1: x is not on any chain. The probability that $f^k(x)$ hits a known element is at most (i - 1)/(N - (i - 1)), because x is uniformly distributed among the N - (i - 1) unknown elements. (Since f is 1-1, it follows that $f^k(x)$ is uniformly distributed over a set of N - (i - 1) elements.)

- Case 2: x is on one chain and $f^k(x)$ hits another chain. The probability to hit an element of another chain (which must belong to the same cycle) is $(i-1)/(N'-(i-1)^2)$, where $N' \ge N/3$ is the number of vertices on the cycle (on which x reside). This is because the chains on the same cycle may be though of having a random relative shift (which ignore the collisions of known vertices). For $i < \sqrt{N/2}$, we obtain a probability bound of $i/\Omega(N)$.
- Case 3: x is on some chain and $f^k(x)$ hits the same chain. Without loss of generality, suppose that $f^k(x) = x$. For this to happen, the length N' of the cycle (on which x reside) must divide k. We upper-bound the probability that all prime factors of N' are prime factors of k. Recall that N' is uniformly selected in [(N/3), (2N/3)], let $P = P_k$ denote the set of prime factors of k, and note that |P| = poly(n) (by the hypothesis $k \in [2^{\text{poly}(n)}]$). We bound the number of integers in [N] that have all prime factors in P by bounding, for every $t \in [n]$, the product of the number of integers in $[2^t]$ with all prime factors in $P' \stackrel{\text{def}}{=} \{p \in P : p < n^c\}$ and the number of (n - t)-bit integers with all prime factors in $P'' \stackrel{\text{def}}{=} P \setminus P'$, where c is a suitable constant (i.e., satisfying $|P| < n^{c-1}$). For $t > n/\log n$, the size of the first set can be upper-bounded by the number of n^c -smooth numbers in $[2^t]^{32}$ which in turn is bounded by $2^{t-(t/c)+o(t)} = 2^{(1-(1/c))\cdot t+o(t)}$. The size of the second set is upper-bounded by $\binom{|P''|}{(n-t)/(c\log n)} < 2^{(1-(1/c))\cdot(n-t)}$, where the inequality uses $|P''| < n^{c-1}$. Thus, we upper-bound the probability that an uniformly chosen integer in [(N/3), (2N/3)] has all prime factors in P by

$$\sum_{t=1}^{n/\log n} 2^{-(1/c) \cdot (n-t)} + \sum_{t=(n/\log n)+1}^{n} 2^{-(1/c) \cdot (t+(n-t)) + o(t)} = 2^{-(n/c) + o(n)}$$

Thus, the probability that we form a collision in q queries (to the implementation) is at most $q^2 \cdot N^{-1/(c+1)}$.

We now turn to the case that f is a random function (as in the specification). Suppose that we make the non-redundent query (x, k). We wish to upper-bound the probability that $f^k(x) = y$, for some fixed y (which is on one of the chains). It is well-known that the expected number of ancestors of y under a random f is $\Theta(\sqrt{N})$; see, e.g., Theorem 33 in [6, Ch. XIV]. Thus, $\mathbf{Pr}_f[|\cup_{i\geq 1} f^{-i}(y)| > N^{3/4}] = O(N^{-1/4})$, and it follows that $\mathbf{Pr}_f[f^k(x) = y] < N^{-1/4} + O(N^{-1/4})$, for any fixed (x, k) and y. (Indeed, it seems that this is a gross over-estimate, but it suffices for our purposes.) It follows that the probability that we form a collision in q queries to the specification is at most $O(q^2/N^{1/4})$.

Comment: The proof of Theorem 9.1 can be easily adapted so to provide a truthful closeimplementation of a random permutation with iterated-evaluation and iterated-inverse queries. That is, we refer to a specifying machine that uniformly selects a permutation $p : \{0,1\}^n \to \{0,1\}^n$, and answers queries of the form (x,m), where $x \in \{0,1\}^n$ and $m \in [\pm 2^{\text{poly}(n)}]$, with the value $p^m(x)$. The implementation is exactly the one used in the proof of Theorem 9.1, and thus we should only analyze the probability of collision when making (non-redundent) queries to a random permutation. For any fixed (x, k) and y, the probability that $\pi^k(x) = y$ equals the probability that x and y resides on the same cycle of the permutation p and that their distance on this cycle equals $k \mod \ell$, where ℓ is the length of this cycle. The claim follows using the fact that ℓ is distributed uniformly over [N]

³²An integer is called *y*-smooth if all its prime factors are smaller that *y*. The fraction of *y*-smooth integers in [x] is upper-bounded by $u^{-u+o(u)}$, where $u = (\log x)/(\log y)$; see, [7]. Thus, in case $t > n/\log n$, the fraction of n^c -smooth integers in $[2^t]$ is upper-bounded by $2^{-(1-o(1))\cdot(t/(\log 2n))\cdot\log_2 t} = 2^{-(1-o(1))t/c}$.

(becuase the probability that x resides on a cycle of a certain length equals the expected number of elements residing on cycles of such length divided by N). An alternative implementation of a random permutation supporting iterated-evaluation (and iterated-inverse) queries was suggested independently by Tsaban [32]. Interestingly, his implementation works by selecting a cycle structure with distribution that is statistically-close to that in a random permutation (and using a set of cycles of corresponding lengths, rather than always using two cycles as we do).

Preimage queries to a random mapping: We turn back to random length preserving functions. Such a random function $f : \{0,1\}^n \to \{0,1\}^n$ is highly unlikely to be 1-1, still the set of preimages of an element under the function is well-defined (i.e., $f^{-1}(y) = \{x : f(x) = y\}$). Indeed, this set may be empty, be a singleton or contain more than one preimage. Furthermore, with overwhelmingly high probability, all these sets are of size at most n. The corresponding "inverse" queries are thus natural to consider.

Theorem 9.2 There exists a truthful close-implementation of the following specification. The specifying machine, uniformly selects a function $f : \{0,1\}^n \to \{0,1\}^n$, and, in addition to the standard evaluation queries, answers the inverse-query $y \in \{0,1\}^n$ with the value $f^{-1}(y)$.

Proof: We start with a truthful implementation that is not statistically-indistinguishable from the specification, but is close to being so and does present our main idea. For $\ell = O(\log n)$ (to be determined), we consider an implementation that uses the orcale in order to define two permutations π_1 and π_2 over $\{0,1\}^n$ (along with their inverses) as well as a random function $g: \{0,1\}^n \to \{0,1\}^\ell$. We define $f(x) = \pi_2(\operatorname{pref}_{n-\ell}(\pi_1(x))g(\pi_1(x)))$, where $\operatorname{pref}_i(z)$ denotes the *i*-bit long prefix of *z*. That is, the function *g* induces collisions within the structured sets S_α , where $S_\alpha \stackrel{\text{def}}{=} \{\alpha\beta : \beta \in \{0,1\}^\ell\}$, and the permutation π_1 (resp., π_2) randomly route inputs (resp., outputs) to (resp., from) these sets. Indeed, it is instructive to note that *g* induces a collection of random independent functions $g_\alpha: \{0,1\}^\ell \to \{0,1\}^\ell$ such that $g_\alpha(\beta) = g(\alpha\beta)$, and that each g_α induces a random function on the corresponding S_α (i.e., mapping $\alpha\beta$ to $\alpha g_\alpha(\beta)$). Thus, letting $\operatorname{suff}_i(z)$ denote the *i*-bit long suffix of *z*, we may write

$$f(x) = \pi_2(\alpha g_\alpha(\beta)), \text{ where } \alpha \leftarrow \operatorname{pref}_{n-\ell}(\pi_1(x)) \text{ and } \beta \leftarrow \operatorname{suff}_{n-\ell}(\pi_1(x)).$$
(3)

The evaluation queries are answered in a straightforward way (i.e., by evaluating π_1 , g and π_2). The inverse-query y is answered by first computing $\alpha\beta = \pi_2^{-1}(y)$, where $|\alpha| = n - \ell$, then computing $R_{\alpha}(\beta) \stackrel{\text{def}}{=} \{\beta' : g(\alpha\beta') = \beta\}$ via exhaustive search, and finally setting $f^{-1}(y) = \{\pi_1^{-1}(\alpha\beta') : \beta' \in R_{\alpha}(\beta)\}$. Indeed, the key point is that, since $\ell = O(\log n)$, we can afford to determine the set $R_{\alpha}(\beta)$ by going over all possible $\beta' \in \{0,1\}^{\ell}$ and including β' if and only if $g(\alpha\beta') = \beta$. The random permutation π_1 (resp., π_2) guarantees that it is unlikely to make two evaluation queries (resp., inverse-queries) that are served via the same set S_{α} . It is also unlikely to have a non-obvious "interaction" between these two types of queries (where an obvious interaction is obtained by asking for a preimage of an answer to an evaluation query or vice versa). Thus, the answers to the evaluation queries look random, and the answers to the inverse-queries are almost independent random subsets with sizes that corresponds to the collision of 2^{ℓ} elements (i.e., 2^{ℓ} balls thrown at random to 2^{ℓ} cells).

The only thing that is wrong with the above implementation is that the sizes of the preimage-sets correspond to the collision pattern of 2^{ℓ} balls thrown at random to 2^{ℓ} cells, rather than to that of the collision pattern of 2^{n} balls thrown at random to 2^{n} cells. Let $p_{i}(m)$ denote the expected fraction

of cells that contain *i* balls, when we throw at random *m* balls into *m* cells. Then, $p_0(m) \approx 1/e$, for all sufficiently large *m*, whereas

$$p_i(m) \approx \frac{1}{(i!)e} \cdot \prod_{j=1}^i \left(1 - \frac{j-2}{m-1}\right)$$
 (4)

We focus on $i \leq n$ (because for i > n both $p_i(2^{\ell})$ and $p_i(2^n)$ are smaller than 2^{-2n}). We may ignore the (negligible in n) dependence of $p_i(2^n)$ on 2^n , but not the (noticeable) dependence of $p_i(2^{\ell})$ on $2^{\ell} = \text{poly}(n)$. Specifically, we have:

i	$p_i(2^n) \sim e^{-1}/(i!)$	$p_i(n^c+1) \approx (\prod_{j=1}^i (1-(j-2)n^{-c})) \cdot p_i(2^n)$ $\approx (\prod_{j=1}^i (1-(j-2)n^{-c})) \cdot (e^{-1}/(i!))$
	$\sim \epsilon /(\iota)$	$\sim (\prod_{j=1}^{n} (1 - (j - 2)\pi)) \cdot (e^{-j} (i:))$
1	e^{-1}	$\frac{(1+n^{-c})\cdot e^{-1}}{1-1}$
2	$e^{-1}/2$	$(1+n^{-c}) \cdot e^{-1}/2$
3	$e^{-1}/6$	$\approx (1 - n^{-2c}) \cdot e^{-1}/6$
4	$e^{-1}/24$	$\approx (1 - 1.5n^{-c}) \cdot e^{-1}/24$
$i \ge 4$	$e^{-1}/(i!)$	$(1 - \Theta(i^2 n^{-c})) \cdot e^{-1}/(i!)$

Thus, the singleton and two-element sets are slightly over-represented in our implementation (when compared to the specification), whereas the larger sets are under-represented. In all cases, the deviation is by a factor related to $1 \pm (1/\text{poly}(n))$, which cannot be tolerated in a close-implementation. Thus, all that is required is to modify the function g such that it is slightly more probable to form larger collisions (inside the sets S_{α} 's). We stress that we can easily compute all the relevant quantities (i.e., all $p_i(2^n)$'s and $p_i(2^{\ell})$'s, for i = 1, ..., n), and so obtaining a close-implementation is merely a question of details, which are shortly outlined next.

Let us just sketch one possible approach. For $N \stackrel{\text{def}}{=} 2^n$ and $t \stackrel{\text{def}}{=} 2^\ell$, we have N/t sets S_{α} 's that are each partitioned at random by the g_{α} 's to subsets (which correspond to the sets of $\alpha\beta$'s that are mapped to the same image under g_{α}). Now, for a random collection of g_{α} 's, the number of *i*-subsets divided by N is $p_i \stackrel{\text{def}}{=} p_i(t)$ rather than $q_i \stackrel{\text{def}}{=} p_i(N)$ as desired. Recall that $|p_i - q_i| \le p_i/(t-1)$ for all $i \geq 1$, and note that $\sum_{i} p_{i}i = 1 = \sum_{i} q_{i}i$. Indeed, it is instructive to consider the fractional mass of elements that resides in *i*-subsets; that is, let $p'_i = p_i i$ and $q'_i = q_i i$. We need to move a fractional mass of about 1/(t-1)e elements from singleton subsets (resp., two-element subsets) to the larger subsets. With overwhelmingly high probability, each S_{α} contains more than n singleton subsets (resp., n/2 two-element subsets). We are going to use only these subsets towards the correction of the distribution of mass; this is more than enough, because we need to relocate only a fractional mass of 1/(t-1)e from each type of subsets (i.e., less than one element per a set S_{α} , which in turn has cardinality t). In particular, we move a fractional mass of $p'_1 - q'_1 = p'_2 - q'_2$ from singleton (resp., two-element) subsets into larger subsets. Specifically, for each $i \ge 3$, we move a fractional mass of $(q'_i - p'_i)/2$ elements residing in singletons and $(q'_i - p'_i)/2$ elements residing in two-element subsets into i-subsets.³³ This (equal contribution condition) will automatically guarantee that the mass in the remaining singleton and two-element subsets is as desired. We stress that there is no need to make the "mass distribution correction process" be "nicely distributed" among the various sets

³³For example, we move mass into 3-subsets by either merging three singletons or merging a singleton and a twosubset into a corresponding 3-subset, where we do three merges of the latter type per each merge of the former type. Similarly, for each $i \ge 4$, we move mass into *i*-subsets by merging either *i* singletons or i/2 two-subsets, while doing an equal number of merges of each type. Finally, for every $j \ge 1$, we move mass into (2j + 3)-subsets by merging additionally created 2j-subsets and 3-subsets (where additional 2-subsets are created by either using a 2-subset or merging two singletons, in equal proportions).

 S_{α} 's, because its affect is anyhow hidden by the application of the random permutation π_2 . The only thing we need is to perform this correction procedure efficiently (i.e., for every α we should efficiently decide how to modify g_{α}), and this is indeed doable.

10 Conclusions and Open Problems

The questions that underlie our work refer to the existence of good implementations of various specifications. At the very least, we require the implementations to be computationally-indistinguishable from the corresponding specifications.³⁴ That is, we are interested in pseudo-implementations. Our ultimate goal is to obtain such implementations via ordinary (probabilistic polynomial-time) machines, and so we ask:

Q1: Which specifications have truthful pseudo-implementations (by ordinary machines)?

- **Q2:** Which specifications have almost-truthful pseudo-implementations (by ordinary machines)?
- Q3: Which specifications have pseudo-implementations at all?

In view of Theorem 2.9, as far as Questions Q1 and Q3 are concerned, we may as well consider implementations by oracle machines (having access to a random oracle). Indeed, the key observation that started us going was that the following questions are the "right" ones to ask:

- Q1r (Q1 revised): Which specifications have truthful close-implementations by oracle machines (having access to a random oracle)?
- Q3r (Q3 revised): Which specifications have such close-implementations at all?

We remark that even in case of Question Q2, it may make sense to study first the existence of implementations by oracle machines, bearing in mind that the latter cannot provide a conclusive positive answer (as shown in Theorem 2.11).

In this work, we have initiated a comprehensive study of the above questions. In particular, we provided a fair number of non-trivial implementations of various specifications relating to the domains of random functions, random graphs and random codes. The challenge of characterizing the class of specifications that have good implementations (e.g., Questions Q1r and Q3r) remains wide open. A good start may be to answer such questions when restricted to interesting classes of specifications (e.g., the class of specifications of random graphs having certain type of properties).

Limited-independence implementations. Our definition of pseudo-implementation is based on the notion of computational indistinguishability (cf. [20, 33, 16]) as a definition of similarity among objects. A different notion of similarity underlies the construction of sample spaces having limited-independence properties (see, e.g., [2, 8]). For example, we say that an implementation is k-wise close to a given specification if the distribution of the answers to any k fixed queries to the implementation is statistically close to the distribution of these answers in the specification. The study of Question Q1r is also relevant to the construction of truthful k-wise close implementations, for any k = poly(n). In particular, one can show that any specification that has a truthful closeimplementation by an oracle machine, has a truthful k-wise close implementation by an ordinary probabilistic polynomial-time machine.³⁵ A concrete example appears at the end of Section 5.

 $^{^{34}}$ Without such a qualification, the questions stated below are either meaningless (i.e., every specification has a "bad" implementation) or miss the point of generating random objects.

³⁵The claim follows by combining an implementation (by an oracle machine) that makes at most t queries to its random oracle with a sample space of $k \cdot t$ -wise independent functions.

Acknowledgments

The first two authors wish to thank Silvio Micali for discussions that took place two decades ago. The main part of Theorem 2.9 was essentially observed in these discussions. These discussions reached a dead-end because the notion of a specification was missing (and so it was not understood that the interesting question is which specifications can be implemented at all (i.e., even by an oracle machine having access to a random function)).

We are grateful to Noga Alon for very helpful discussions regarding random graphs and explicit constructions of bounded-degree graphs of logarithmic girth. We also thank Avi Wigderson for a helpful discussion regarding the proof of Lemma 6.3. Finally, thanks to Moni Naor for calling our attention to [11], and to Omer Reingold and S. Muthu for calling our attention to [13, Lem. 2].

References

- M. Abadi, E. Allender, A. Broder, J. Feigenbaum, and L. Hemachandra. On Generating Hard, Solved Instances of Computational Problem. In *Crypto88*, pages 297–310.
- [2] N. Alon, L. Babai and A. Itai. A fast and Simple Randomized Algorithm for the Maximal Independent Set Problem. J. of Algorithms, Vol. 7, pages 567–583, 1986.
- [3] E. Bach. Analytic Methods in the Analysis and Design of Number-Theoretic Algorithms. ACM Distinguished Dissertation (1984), MIT Press, Cambridge MA, 1985.
- [4] S. Ben-David, B. Chor, O. Goldreich, and M. Luby. On the Theory of Average Case Complexity. JCSS, Vol. 44, No. 2, 1992, pages 193–219. Preliminary version in 21st STOC, 1989.
- [5] M. Blum and S. Micali. How to Generate Cryptographically Strong Sequences of Pseudo-Random Bits. SICOMP, Vol. 13, pages 850–864, 1984. Preliminary version in 23rd FOCS, 1982.
- [6] B. Bollobas. Random Graphs. Academic Press, 1985.
- [7] E.R. Canfield, P. Erdos, and C. Pomerance. On a Problem of Oppenheim Concerning "Factorisatio Numerorum". Jour. of Number Theory, Vol. 17, pages 1–28, 1983.
- [8] B. Chor and O. Goldreich. On the Power of Two-Point Based Sampling. Jour. of Complexity, Vol 5, 1989, pages 96-106. Preliminary version dates 1985.
- [9] I. Damgard. Collision Free Hash Functions and Public Key Signature Schemes. In Euro-Crypt'87, Springer-Verlag, LNCS 304, pages 203-216.
- [10] P. Erdos and H. Sachs. Reguläre Graphen gegenebener Taillenweite mit minimaler Knotenzahl. Wiss. Z. Univ. Halle-Wittenberg, Math. Nat. R., 12, pages 251–258, 1963.
- [11] J. Feigenbaum, S. Kannan, M. Strauss, M. Viswanathan. An Approximate L1-Difference Algorithm for Massive Data Streams. Proceedings of 40th FOCS, pages 501-511, 1999.
- [12] P. Flajolet and A.M. Odlyzko. Random mapping statistics. In *EuroCrypt'89*, Springer-Verlag, LNCS 434, pages 329–354.
- [13] A. Gilbert, S. Guha, P. Indyk, Y. Kotidis, S. Muthukrishnan, and M. Strauss. Fast, Small-Space Algorithms for Approximate Histogram Maintenance. In the proceedings of 34th STOC, pages 389–398, 2002.
- [14] O. Goldreich. A Note on Computational Indistinguishability. *IPL*, Vol. 34, pages 277–281, May 1990.
- [15] O. Goldreich. Foundation of Cryptography Basic Tools. Cambridge University Press, 2001.
- [16] O. Goldreich, S. Goldwasser, and S. Micali. How to Construct Random Functions. JACM, Vol. 33, No. 4, pages 792–807, 1986.
- [17] O. Goldreich, and H. Krawczyk, On Sparse Pseudorandom Ensembles. Random Structures and Algorithms, Vol. 3, No. 2, (1992), pages 163–174.

- [18] O. Goldreich and D. Ron. Property Testing in Bounded Degree Graphs. Algorithmica, 32 (2), pages 302-343, 2002.
- [19] O. Goldreich and L. Trevisan. Three Theorems regarding Testing Graph Properties. Proceedings of 42nd FOCS, pages 460-469, 2001. Full version in ECCC, TR01-010, 2001.
- [20] S. Goldwasser and S. Micali. Probabilistic Encryption. JCSS, Vol. 28, No. 2, pages 270–299, 1984. Preliminary version in 14th STOC, 1982.
- [21] J. Hastad, R. Impagliazzo, L.A. Levin and M. Luby. A Pseudorandom Generator from any One-way Function. SICOMP, Volume 28, Number 4, pages 1364–1396, 1999. Preliminary versions by Impagliazzo et. al. in 21st STOC (1989) and Hastad in 22nd STOC (1990).
- [22] R. Impagliazzo and A. Wigderson. P=BPP if E requires exponential circuits: Derandomizing the XOR Lemma. In 29th STOC, pages 220–229, 1997.
- [23] W. Imrich. Explicit Construction of Regular Graphs with no Small Cycles. Combinatorica, Vol. 4, pages 53–59, 1984.
- [24] S. Janson. The numbers of spanning trees, Hamilton cycles and perfect matchings in a random graph. Combin. Prob. Comput., Vol. 3, pages 97–126, 1994.
- [25] D.E. Knuth. The Art of Computer Programming, Vol. 2 (Seminumerical Algorithms). Addison-Wesley Publishing Company, Inc., 1969 (first edition) and 1981 (second edition).
- [26] F. Lazebnik and V.A. Ustimenko. Explicit Construction of Graphs with arbitrary large Girth and of Large Size.
- [27] L.A. Levin. Average Case Complete Problems. SICOMP, Vol. 15, pages 285–286, 1986.
- [28] A. Lubotzky, R. Phillips, P. Sarnak, Ramanujan Graphs. Combinatorica, Vol. 8, pages 261–277, 1988.
- [29] M. Luby and C. Rackoff. How to Construct Pseudorandom Permutations from Pseudorandom Functions. SICOMP, Vol. 17, 1988, pages 373–386.
- [30] G.A. Margulis. Explicit Construction of Graphs without Short Cycles and Low Density Codes. Combinatorica, Vol. 2, pages 71–78, 1982.
- [31] M. Naor and O. Reingold. Constructing Pseudo-Random Permutations with a Prescribed Structure, Jour. of Crypto., Vol. 15 (2), 2002, pages 97–102.
- [32] B. Tsaban. Permutation graphs, fast forward permutations, and sampling the cycle structure of a permutation. *Journal of Algorithms*, Vol. 47 (2), pages 104–121, 2003.
- [33] A.C. Yao. Theory and Application of Trapdoor Functions. In 23rd FOCS, pages 80-91, 1982.

Appendix A: Implementing various probability distributions

Our proof of Theorem 5.2 relies on efficient procedures for generating elements from a finite set according to two probability distributions. In both cases, we need procedures that work in time that is poly-logarithmic (rather than polynomial) in the size of the set (and the reciprocal of the desired approximation parameter). In both cases, we have close expressions (which can be evaluated in poly-logarithmic time) for the probability mass that is to be assigned to each element. Thus, in both cases, it is easy to generate the desired distribution in time that is almost-linear in the size of the set. Our focus is on generating good approximations of these distributions in time that is poly-logarithmic in the size of the set.

Indeed, the problem considered in this appendix is a special case of our general framework. We are given a specification of a distribution (i.e., each query should be answered by a sample drawn independently from that distribution), and we wish to closely-implement it (i.e., answer each query by a sample drawn independently from approximately that distribution).

A.1 Sampling the binomial distribution

We first consider the generation of elements according to the binomial distribution. For any N, we need to output any value $v \in \{0, 1, ..., N\}$ with probability $\binom{N}{v} \cdot 2^{-N}$. An efficient procedure for this purpose is described in Knuth [25, Sec. 3.4.1]. In fact, Knuth describes a more general procedure that, for every p, outputs the value $v \in \{0, 1, ..., N\}$ with probability $b_{N,p}(v) \stackrel{\text{def}}{=} \binom{N}{v} \cdot p^v (1-p)^{N-v}$. However, his description is in terms of operations with reals, and so we need to adapt it to the standard (bit-operation) model. Knuth's description proceeds in two steps:

- 1. In Section 3.4.1.F, it is shown how to reduce the generation of the binomial distribution $b_{N,p}$ to the generation of some *beta distributions*, which are continuous distributions over [0, 1] that depends on two parameters a and b.³⁶ The reduction involves taking $\log_2 N$ samples from certain beta distributions, where the parameters of these distributions are easily determined as a function of N. The samples of the beta distributions are processed in a simple manner involving only comparisons and basic arithmetic operations (subtraction and division).
- 2. In Section 3.4.1.E, it is shown how to generate any beta distribution. The generator takes a constant number of samples from the continuous uniform distribution over [0, 1], and produces the desired sample with constant probability (otherwise, the process is repeated). The samples of the uniform distributions are processed in a simple manner involving only comparisons and various arithmetic and trigonometric operations (including computing functions as log and tan).

The above is described in terms of real arithmetic and sampling uniformly in [0, 1], and provides a perfect implementation. The question is what happens when we replace the samples with ones taken from the set $\{\epsilon, 2\epsilon, ..., \lfloor 1/\epsilon \rfloor \cdot \epsilon\}$, and replace the real arithmetics with approximations up to a factor of $1 \pm \epsilon$.

$$F_{a,b}(r) \stackrel{\text{def}}{=} a \cdot \binom{a+b-1}{a} \cdot \int_0^r x^{a-1} (1-x)^{b-1} dx$$

 $^{^{36}}$ A beta distribution with (natural) parameters a and b is defined in terms of the accumulative distribution function

and the uniform continuous distribution is a special case (i.e., a = b = 1). In general, $F_{a,b}(r)$ equals the probability that the *b*th largest of a + b - 1 independent uniformly chosen samples in [0, 1] has value at most r.

Let us first consider the effect of replacing the uniform continuous distribution U(r) = r by the continuous step-distribution $S_{\epsilon}(r) \stackrel{\text{def}}{=} \lfloor r/\epsilon \rfloor \cdot \epsilon$, where we may assume that $1/\epsilon$ is an integer. Since the variation distance between U and S_{ϵ} is $O(\epsilon)$, the same holds for any function applied to a constant number of samples taken from these distribution. Thus, the implementation of the beta distributions via the step-distribution S_{ϵ} will deviate by only $O(\epsilon)$, and using the latter to generate the binomial distribution $b_{N,p}$ only yields a deviation of $O(\epsilon \log N)$. Finally, using the *average* numerical stability of all functions employed³⁷ we conclude that an implementation by $O(\log(1/\epsilon))$ bits of precision will only introduce a deviation of ϵ .

A.2 Sampling from the two-set total-sum distribution

We now turn to the generation of pairs (l, r) such that l + r = T and $0 \le l, r \le S$, where $T \le 2S$. Specifically, we need to produce such a pair with probability proportional to $\binom{S}{l} \cdot \binom{S}{r}$ (i.e., the number of ways to select l elements from one set of size S and r elements from another such set). (In the proof of Theorem 5.2, S = M/2.) Without loss of generality, we may assume that $T \le S$ (or else we select the "complementary" elements). Thus, we need to sample $r \in \{0, ..., T\}$ with probability

$$p_r = \frac{\binom{S}{T-r} \cdot \binom{S}{r}}{\binom{2S}{T}}$$
(5)

We wish to produce a sample with deviation at most ϵ from the correct distribution and are allowed time poly(k), where $k \stackrel{\text{def}}{=} \log(S/\epsilon)$. In case $T \leq k$, we perform this task in the straightforward manner; that it, compute all the T + 1 probabilities p_r , and select r accordingly. Otherwise (i.e., T > k), we rely on the fact that p_r is upper-bounded by twice the binomial distribution of Ttries (i.e., $q_r = {T \choose r}/2^T$). This leads to the following sampling process:

- 1. Select r according to the binomial distribution of T tries.
- 2. Compute p_r and q_r . Output r with probability $p_r/2q_r$, and go to Step 1 otherwise.

We will show (see Fact A.1 below) that $p_r \leq 2q_r$ always holds. Thus, in each iteration, we output r with probability that is proportional to p_r ; that is, we output r with probability $q_r \cdot (p_r/2q_r) = p_r/2$. It follows that each iteration of the above procedure produces an output with probability 1/2, and by truncating the procedure after k iterations (and producing arbitrary output in such a case) the output distribution is statistically close to the desired one.

Fact A.1 Suppose that $T \leq S$ and T > k. For p_r 's and q_r 's as above, it holds that $p_r < 2q_r$.

Proof: The cases r = T and r = 0 are readily verified (by noting that $p_r = {S \choose T} / {2S \choose T} < 2^{-T}$ and $q_r = 2^{-T}$). For $r \in \{1, ..., T-1\}$, letting $\alpha \stackrel{\text{def}}{=} (S-r)/(2S-T) \in (0,1)$, we have

$$\frac{p_r}{q_r} = \frac{\binom{S}{r} \cdot \binom{S}{T-r} / \binom{2S}{T}}{\binom{T}{r} / 2^T} = 2^T \cdot \frac{\binom{2S-T}{S-r}}{\binom{2S}{S}}$$

³⁷Each of these functions (i.e., rational expressions, log and tan) has a few points of instability, but we apply these functions on arguments taken from either the uniform distribution or the result of prior functions on that distribution. In particular, except for what happens in an ϵ -neighborhood of some problematic points, all functions can be well-approximated when their argument is given with $O(\log(1/\epsilon))$ bits of precision. Furthermore, the functions log and tan are only evaluated at the uniform distribution (or simple functions of it), and the rational expressions are evaluated on some intermediate beta distributions. Thus, in all cases, the problematic neighborhoods are only assigned small probability mass (e.g., ϵ in the former case and $O(\sqrt{\epsilon})$ in the latter).

$$= 2^{T} \cdot (1+o(1)) \cdot \frac{(2\pi\alpha(1-\alpha)\cdot(2S-T))^{-1/2}\cdot 2^{\mathbf{H}_{2}(\alpha)\cdot(2S-T)}}{(2\pi(1/2)^{2}\cdot 2S)^{-1/2}\cdot 2^{\mathbf{H}_{2}(1/2)\cdot 2S}}$$
$$= \frac{1+o(1)}{\sqrt{2\alpha(1-\alpha)\cdot\beta}} \cdot 2^{(\mathbf{H}_{2}(\alpha)-1)\cdot(2S-T)}$$

where $\beta \stackrel{\text{def}}{=} (2S - T)/S \ge 1$ and \mathbf{H}_2 is the binary entropy function. For $\alpha \in [(1/3), (2/3)]$, we can upper-bound p_r/q_r by $(1 + o(1)) \cdot \sqrt{9/4\beta} < 2$. Otherwise (i.e., without loss of generality $\alpha < 1/3$), we get that $\mathbf{H}_2(\alpha) < 0.92$ and $\alpha^{-1}(1 - \alpha)^{-1} \le 2S - T$, where for the latter inequality we use $1 \le r \le S - 1$. Thus, p_r/q_r is upper-bounded by $O(\sqrt{2S - T}) \cdot 2^{-\Omega(2S - T)} = O(2^{-\Omega(S) + \log S})$, which vanishes to zero with k (because $S \ge T > k$).³⁸

A.3 A general tool for sampling strange distributions

In continuation to Appendix A.2, we state a useful lemma (which was implicitly used above as well as in prior works). The lemma suggests that $poly(\log N)$ -time sampling from a desired probability distribution $\{p_i\}_{i=1}^N$ can be reduced to sampling from a related probability distribution $\{q_i\}_{i=1}^N$, which is hopefully $poly(\log N)$ -time sampleable.

Lemma A.2 Let $\{p_i\}_{i=1}^N$ and $\{q_i\}_{i=1}^N$ be probability distributions satisfying the following conditions:

- 1. There exists a polynomial-time algorithm that given $i \in [N]$ outputs approximations of p_i and q_i up to $\pm N^{-2}$.
- 2. Generating an index i according to the distribution $\{q_i\}_{i=1}^N$ is closely-implementable (upto negligible in log N deviation and in poly(log N)-time).
- 3. There exist a poly(log N)-time recognizable set $S \subseteq [N]$ such that
 - (a) $1 \sum_{i \in S} p_i$ is negligible in $\log N$.
 - (b) There exists a polynomial p such that for every $i \in S$ it holds that $p_i \leq p(\log N) \cdot q_i$.

Then generating an index i according to the distribution $\{p_i\}_{i=1}^N$ is closely-implementable.

Proof: Without loss of generality, S may exclude all i's such that $p_i < N^{-2}$. For simplicity, we assume below that given i we can exactly compute p_i and q_i (rather than only approximate them within $\pm N^{-2}$). Let $t \stackrel{\text{def}}{=} p(\log N)$. The sampling procedure proceeds in iterations, where in each iteration i is selected according to the distribution $\{q_i\}_{i=1}^N$, and is output with probability p_i/tq_i if $i \in S$. (Otherwise, we proceed to the next iteration.) Observe that, conditioned on producing an output, the output of each iteration is in S and equals i with probability $q_i \cdot (p_i/tq_i) = p_i/t$. Thus, each iteration produces output with probability $\sum_{i \in S} p_i/t > 1/2t$, and so halting after $O(t \log(1/\epsilon))$ iterations we produce output with probability at least $1 - \epsilon$. For any $i \in S$, the output is i with probability $(1 \pm \epsilon) \cdot p_i/\rho$, where $\rho \stackrel{\text{def}}{=} \sum_{j \in S} p_j$. Setting ϵ to be negligible in $\log N$, the lemma follows.

A typical application of Lemma A.2 is to the case that for each $i \in [N]$ the value of p_i can be approximated by one out of $m = poly(\log N)$ predetermined p_j 's. Specifically:

³⁸In fact, it holds that $p_r \leq \sqrt{2} \cdot q_r$ for all r's, with the extreme value obtained at r = T/2 (and T = S), where we have $\alpha = 1/2$ (and $\beta = 1$).

Corollary A.3 Let $\{p_i\}_{i=1}^N$ be a probability distribution and $S \subseteq [N]$ be a set satisfying Conditions (1) and (3a) of Lemma A.2. Suppose that, for $m, t = \text{poly}(\log N)$, there exists an efficiently constructible sequence of integers $1 = i_1 < i_2 < \cdots < i_m = N$ such that for every $j \in [m-1]$ and $i \in [i_j, i_{j+1}] \cap S$ it holds that $p_{i_j}/t < p_i < t \cdot p_{i_j}$. Then generating an index *i* according to the distribution $\{p_i\}_{i=1}^N$ is closely-implementable.

Proof: For every $j \in [m-1]$ and $i \in [i_j, i_{j+1}] \cap S$, define $p'_i = p_{i_j}$ and note that $p'_i/t < p_i < t \cdot p'_i$. Let $p' = \sum_{i \in S} p'_i$, and note that p' < t. Now, define $q_i = p'_i/p'$ for every $i \in S$, and $q_i = 0$ otherwise. Then, for every $i \in S$, it holds that $p_i < t \cdot p'_i = t \cdot p' \cdot q_i < t^2 q_i$. Since these q_i 's satisfy Conditions (1), (2) and (3b) of Lemma A.2, the corollary follows.

Appendix B: Implementing a Random Bipartite Graph

Following the description in Section 6, we present a close-implementation of random bipartite graphs. Two issues arise. Firstly, we have to select the proportion of the sizes of the two parts, while noticing that different proportions give rise to different number of graphs. Secondly, we note that a bipartite graph uniquely defines a 2-partition (up to switching the two parts) only if it is connected. However, since all but a negligible fraction of the bipartite graphs are connected, we may ignore the second issue, and focus on the first one. (Indeed, the rest of the discussion is slightly imprecise because the second issue is ignored.)

For $i \in [\pm N]$, the number of 2N-vertex bipartite graphs with N + i vertices on the first part is

$$\binom{2N}{N+i} \cdot 2^{(N+i)\cdot(N-i)} \leq \binom{2N}{N} \cdot 2^{N^2 - i^2}$$

where equality holds for i = 0 and approximately holds (i.e., upto a constant factor) for $|i| = \sqrt{N}$. Thus, all but a negligible fraction of the 2*N*-vertex bipartite graphs have $N \pm \log_2 N$ vertices on each part. That is, we may focus on $O(\log N)$ values of *i*. Indeed, for each $i \in [\pm \log_2 N]$, we compute $T_i \stackrel{\text{def}}{=} \binom{2N}{N+i} \cdot 2^{N^2-i^2}$, and $p_i = T_i/T$, where $T \stackrel{\text{def}}{=} \sum_{j=-\log_2 N}^{\log_2 N} T_j$. Next, we select *i* with probability p_i , and construct a random 2*N*-vertex bipartite graph with N + i vertices on the first part as follows:

- As in Section 6, we use the function f_1 to implement a permutation π . We let $S \stackrel{\text{def}}{=} \{v : \pi(v) \in [N+i]\}$, and $\chi_S(i) \stackrel{\text{def}}{=} 1$ if and only if $i \in S$.
- As in Section 6, we answer the query (u, v) by 0 if $\chi_S(u) = \chi_S(v)$ and according to the value of f_2 otherwise.

Appendix C: Various Calculations

For the proof of Lemma 6.3

The proof of Lemma 6.3 refers to the following known fact:

Fact C.1 Let X be a random variable ranging over some domain D, and suppose that $\mathbf{H}(X) \geq \log_2 |D| - \epsilon$. Then X is at statistical distance at most $O(\sqrt{\epsilon})$ from the uniform distribution over D.

Proof: Suppose that X is at statistical distance δ from the uniform distribution over D. Then, there exists a $S \subset D$ such that $|\mathbf{Pr}[X \in S] - (|S|/|D|)| = \delta$, and assume without loss of generality that $|S| \ge |D|/2$. Note that either for each $e \in S$ it holds that $\mathbf{Pr}[X = e] \ge 1/|D|$ or for each $e \in S$ it holds that $\mathbf{Pr}[X = e] \le 1/|D|$ or for each $e \in S$ it holds that $\mathbf{Pr}[X = e] \le 1/|D|$. By removing the |S| - (|D|/2) elements of smallest absolute difference (i.e., smallest $|\mathbf{Pr}[X = e] - (1/|D|)|$), we obtain a set S' of size |D|/2 such that $|\mathbf{Pr}[X \in S'] - (|S'|/|D|)| \ge \delta/2$. The entropy of X is maximized when it is uniform both on S' and on $D \setminus S'$. Thus:

$$\begin{aligned} \mathbf{H}(X) &\leq \mathbf{H}_2(\mathbf{Pr}[X \in S']) + \mathbf{Pr}[X \in S'] \cdot \mathbf{H}(X|X \in S') + \mathbf{Pr}[X \in D \setminus S'] \cdot \mathbf{H}(X|X \in D \setminus S') \\ &= \mathbf{H}_2\left(\frac{1}{2} + \frac{\delta}{2}\right) + \log_2(|D|/2) \\ &= 1 - \Omega(\delta^2) + \log_2(|D|/2) \end{aligned}$$

We get that $\mathbf{H}(X) \leq \log_2 |D| - c \cdot \delta^2$, for some universal c > 0. Combining this with the hypothesis that $\mathbf{H}(X) \geq \log_2 |D| - \epsilon$, we get that $\epsilon \geq c \cdot \delta^2$, and $\delta \leq \sqrt{\epsilon/c}$ follows.

For the proof of Theorem 6.5

In continuation to Footnote 23, which refers to Part 2 of the proof of Theorem 6.5, we prove the following fact.

Fact C.2 Let c(N) be as in Theorem 6.5, and $T \stackrel{\text{def}}{=} N/c(N)$. Assume that T is an integer. Consider any fixed partition, $(P_1, ..., P_T)$, of [N] such that $|P_i| = c(N)$ for every *i*. Consider a graph selected as follows:

- Each P_i is an independent set.
- For $k = \binom{c(N)+3}{2}$, the rest of the edges are determined by a k-wise independent binary sequence of length $\binom{N}{2} T \cdot \binom{c(N)}{2}$.

Then, with probability at least $1 - O(N^{-1/2})$, the graph has no independent set of size c(N) + 3.

Proof: We will show that the expected number of independet set of size c(N) + 3 is $O(N^{-1/2})$, and the fact will follow. Let $c \stackrel{\text{def}}{=} c(N)$ and $s \stackrel{\text{def}}{=} c + 3$. We partition all possible independent sets of size s into classes according to the contributions of the various P_i 's to them. That is, the classes that corresponds to the sequence $(s_1, ..., s_T)$, where $\sum_{i=1}^T s_i = s$, consists of independent sets having s_i vertices from P_i . For such a class, we let r_j denote the j-th non-zero s_i . We actually, cluster the classes according to the resulting sequence of r_j 's. That is, the cluster $(r_1, ..., r_t)$, where $\sum_{j=1}^t r_j = s$ and $r_j \ge 1$, consists of independent sets having r_j vertices from the j-th part that contributes any vertices to the independent set. Then, the contribution of such a cluster to the expectation is given by the number of potential independent sets in the cluster times the probability that such a potential independent set is assigned no edges. Observe that the number of potential undetermined edges in such a potential independent set is $\binom{s}{2} - \sum_i \binom{r_i}{2}$, and thus the contribution of the cluster is given by

$$\begin{pmatrix} T \\ t \end{pmatrix} \cdot \left[\prod_{i=1}^{t} \binom{c}{r_i} \right] \cdot 2^{-\binom{s}{2} - \sum_{i=1}^{t} \binom{r_i}{2}} = 2^{-\binom{s}{2}} \cdot \binom{N/c}{t} \cdot \prod_{i=1}^{t} \left[\binom{c}{r_i} \cdot 2^{\binom{r_i}{2}} \right]$$

$$< 2^{-\binom{s}{2}} \cdot (N/c)^t \cdot \prod_{i=1}^{t} \left[\binom{c}{r_i} \cdot 2^{\binom{r_i}{2}} \right]$$

We bound, separately, each factor of the form $\binom{c}{r_i} \cdot 2^{\binom{r_i}{2}}$. Specifically:

Claim: Let $f(x) = {\binom{c}{x}} \cdot 2^{\binom{x}{2}}$. Then, for $x \in \{1, ..., c\}$ it holds that $f(x) \leq cN^{x-1}$, and for $x \in \{2, ..., c-1\}$ it holds that $f(x) < N^{x-(3/2)}$.

Using this claim, the contribution of each sequence of r_i 's is:

$$2^{-\binom{s}{2}} \cdot (N/c)^{t} \cdot \prod_{i=1}^{t} f(r_{i}) < 2^{-\binom{s}{2}} \cdot (N/c)^{t} \cdot \prod_{i=1}^{t} cN^{r_{i}-1}$$
$$= 2^{-\binom{s}{2}} \cdot N^{\sum_{i=1}^{t} r_{i}}$$
$$= 2^{-\frac{s(s-1)}{2} + \frac{c}{2} \cdot s} = 2^{-s}$$

where the last two equalities use $\log_2 N = c/2$ and c = s - 3, respectively. Furthermore, if the sequence of r_i 's has some element in $\{2, ..., c-1\}$ then we get a better bound of $2^{-s} \cdot N^{-1/2}$, because we gain at least a factor of $N^{-1/2}$ in the inequality.

Now, the number of sequences $(r_1, ..., r_t)$, for various t, is $\sum_{t=2}^{s} {s-1 \choose t-1} < 2^{s-1}$. Of these only a constant number have all r_i 's in $\{1, c\}$ (i.e., the all-1 sequence and the permutations of (c, 1, 1, 1)). Thus, the expectation is bounded by

$$O(1) \cdot 2^{-s} + 2^s \cdot (2^{-s} \cdot N^{-1/2}) < O(N^{-1/2})$$

where the inequality uses $s > c = 2 \log_2 N$.

We now turn to the proof of the claim. For x = 1, equality holds (i.e., $f(1) = c \cdot 2^0 = c \cdot N^0$). (In fact, this is the only case where equality holds.) For x = c, we have $f(c) = 1 \cdot 2^{c(c-1)/2} = N^{c-1} < cN^{c-1}$. In all other cases, we define $g(x) = \log_2 f(x) - (x - (3/2)) \cdot \log_2 N$, and prove that it is negative. Using $\log_2 N = c/2$, note that

$$g(x) = \log_2 {\binom{c}{x}} + \frac{(x-1)x}{2} - (x-(3/2)) \cdot \frac{c}{2}$$

= $\log_2 {\binom{c}{x}} + \frac{(x-c-1)x}{2} + \frac{3c}{4}$
< $\log_2 {\binom{c'}{x}} + \frac{(x-c')x}{2} + \frac{3c'}{4}$

where $c' \stackrel{\text{def}}{=} c + 1$. Using the fact that $c = c(N) = \omega(1)$ (and $2 \le x \le c' - 2$), we consider two cases:

1. If either $2 \le x \le 5$ or $c' - 5 \le x \le c' - 2$ then we bound $\log_2 {c' \choose x}$ by $5 \log_2 c'$ and get

$$g(x) < 5 \log_2 c' + \frac{1}{2} \cdot \max_{2 \le x \le c'-2} \{ (x - c')x \} + \frac{3c'}{4}$$

= $5 \log_2 c' - \frac{2(c'-2)}{2} + \frac{3c'}{4} < 0$

2. If $5 \le x \le c' - 5$ then we bound $\log_2 {c' \choose x}$ by c' and get

$$g(x) < c' + \frac{1}{2} \cdot \max_{5 \le x \le c' - 5} \{ (x - c')x \} + c'$$

= $2c' - \frac{5(c' - 5)}{2} < 0$

So the claim follows, and so does the entire fact.

Appendix D: A strengthening of Proposition 2.15

The hypothesis of Part 2 of Proposition 2.15 requires the existence of one-way functions, or equivalently the ability to generate hard-instances (to NP-problems) along with corresponding solutions (cf. [15, Sec 2.1]). A seemingly weaker condition, which is in the spirit of Levin's theory of average-case complexity [27] (see also [4]), is the ability to generate hard-instances to NP-problems. Specifically:

Definition D.1 (generating hard instances): A probabilistic polynomial-time algorithm G is called a generator of hard instances for a set S if for every probabilistic polynomial-time algorithm A the probability that A correctly decides whether or not $G(1^n)$ is in S is bounded away from 1. That is, there exists a polynomial p such that for all sufficiently large n's it holds that

$$\mathbf{Pr}_{x \leftarrow G(1^n)}[A(x) = \chi_S(x)] < 1 - \frac{1}{p(n)}$$

where $\chi_S(x) = 1$ if $x \in S$ and $\chi_S(x) = 0$ otherwise.

Definition D.1 only requires that hard instances be generated with "noticible" probability. Note that the existence of one-way functions (even weak ones) implies the ability to generate hard instances to NP-problems. The converse is not known to hold. Thus, the following result strengthens Part 2 of Proposition 2.15.

Proposition D.2 Assuming the existence of generators of hard instances for NP-problems, there exist specifications that cannot be pseudo-implemented.

Proof: Let L be an NP-set that has a generator G of hard instances, let R be the corresponding witness relation (i.e., $L = \{x : \exists y \text{ s.t. } (x, y) \in R\}$), and $R(x) \stackrel{\text{def}}{=} \{y : (x, y) \in R\}$. Consider the specification that answers query x with a uniformly distributed $y \in R(x)$ if $R(x) \neq \emptyset$ and with a special symbol otherwise. We will show that this specification cannot be pseudo-implemented.

Let I be an arbitrary implementation of the above specification, and consider a distinguisher that, for parameter n, makes the query $x \leftarrow G(1^n)$, obtains the answer y, and outputs 1 if and only if $(x, y) \in R$ (which is polynomial-time decidable). When this distinguisher queries the specification, it outputs 1 with probability that equals $\rho \stackrel{\text{def}}{=} \mathbf{Pr}[G(1^n) \in L]$. Assume, towards the contradiction, that when the distinguisher queries I it outputs 1 with probability that at least $\rho - \mu(n)$, where μ is a negligible function. In such a case we obtain a probabilistic polynomial-time algorithm that violates the hypothesis: Specifically, consider an algorithm A such that A(x) answers 1 if and only if $(x, I(x)) \in R$, and note that A is always correct when it outputs 1. Thus,

$$\begin{aligned} \mathbf{Pr}_{x \leftarrow G(1^n)}[A(x) &= \chi_L(x)] &= \mathbf{Pr}[x \in L \land A(x) = 1] + \mathbf{Pr}[x \notin L] \cdot \mathbf{Pr}[A(x) = 0 | x \notin L] \\ &= \mathbf{Pr}[x \in L \land (x, I(x)) \in R] + (1 - \rho) \cdot \mathbf{Pr}[(x, I(x)) \notin R | x \notin L] \\ &\geq (\rho - \mu(n)) + (1 - \rho) \cdot 1 = 1 - \mu(n) \end{aligned}$$

Thus, the implementation I cannot be computationally indistinguishable from the specification, and the proposition follows.