

An Efficient Protocol for Secure Two-Party Computation in the Presence of Malicious Adversaries

Yehuda Lindell^{1,*} and Benny Pinkas^{2,**}

¹ Dept. of Computer Science, Bar-Ilan University, Israel
lindell@cs.biu.ac.il

² Dept. of Computer Science, University of Haifa, Israel
benny@pinkas.net

Abstract. We show an efficient secure two-party protocol, based on Yao's construction, which provides security against malicious adversaries. Yao's original protocol is only secure in the presence of semi-honest adversaries. Security against malicious adversaries can be obtained by applying the compiler of Goldreich, Micali and Wigderson (the "GMW compiler"). However, this approach does not seem to be very practical as it requires using generic zero-knowledge proofs.

Our construction is based on applying cut-and-choose techniques to the original circuit and inputs. Security is proved according to the *ideal/real simulation paradigm*, and the proof is in the standard model (with no random oracle model or common reference string assumptions). The resulting protocol is computationally efficient: the only usage of asymmetric cryptography is for running $O(1)$ oblivious transfers for each input bit (or for each bit of a statistical security parameter, whichever is larger). Our protocol combines techniques from folklore (like cut-and-choose) along with new techniques for efficiently proving consistency of inputs. We remark that a naive implementation of the cut-and-choose technique with Yao's protocol does *not* yield a secure protocol. This is the first paper to show how to properly implement these techniques, and to provide a full proof of security.

Our protocol can also be interpreted as a constant-round black-box reduction of secure two-party computation to oblivious transfer and perfectly-hiding commitments, or a black-box reduction of secure two-party computation to oblivious transfer alone, with a number of rounds which is linear in a statistical security parameter. These two reductions are comparable to Kilian's reduction, which uses OT alone but incurs a number of rounds which is linear in the depth of the circuit [18].

1 Introduction

Secure two-party computation. In the setting of two-party computation, two parties with respective private inputs x and y , wish to jointly compute a functionality

* Research supported in part by an Infrastructures grant from the Ministry of Science, Israel.

** Research supported in part by the Israel Science Foundation (grant number 860/06).

$f(x, y) = (f_1(x, y), f_2(x, y))$, such that the first party receives $f_1(x, y)$ and the second party receives $f_2(x, y)$. Loosely speaking, the security requirements are that nothing is learned from the protocol other than the output (*privacy*), and that the output is distributed according to the prescribed functionality (*correctness*). The actual definition follows the simulation paradigm and blends the above two requirements. Of course, security must be guaranteed even when one of the parties is adversarial. Such an adversary may be *semi-honest* (or passive), in which case it correctly follows the protocol specification, yet attempts to learn additional information by analyzing the transcript of messages received during the execution. In contrast, the adversary may be *malicious* (or active), in which case it can arbitrarily deviate from the protocol specification.

The first general solutions for the problem of secure computation were presented by Yao [29] for the two-party case (with security against semi-honest adversaries) and Goldreich, Micali and Wigderson [11] for the multi-party case (with security even against malicious adversaries). Thus, the results of [29] and [11] constitute important and powerful feasibility results for secure two-party and multi-party computation.

Yao's protocol. In [29], Yao presented a *constant-round* protocol for securely computing any functionality in the presence of semi-honest adversaries. Denote party P_1 and P_2 's respective inputs by x and y and let f be the functionality that they wish to compute (for simplicity, assume that both parties wish to receive $f(x, y)$). Loosely speaking, Yao's protocol works by having one of the parties (say party P_1) first generate a "garbled" (or encrypted) circuit computing $f(x, \cdot)$ and then send it to P_2 . The circuit is such that it reveals nothing in its encrypted form and therefore P_2 learns nothing from this stage. However, P_2 can obtain the output $f(x, y)$ by "decrypting" the circuit. In order to ensure that P_2 learns nothing more than the output itself, this decryption must be "partial" and must reveal $f(x, y)$ only. Without going into unnecessary details, this is accomplished by P_2 obtaining a series of keys corresponding to its input y , such that given these keys and the circuit, the output value $f(x, y)$, and only this value, may be obtained. Of course, P_2 must somehow receive these keys without revealing anything about y to P_1 . This can be accomplished by running $|y|$ instances of a secure 1-out-of-2 Oblivious Transfer protocol [27,7]. Yao's generic protocol is highly efficient, and even practical, for functionalities that have relatively small circuits. An actual implementation of the protocol was presented in [21], with very reasonable performance.

Security against malicious behavior. Yao's protocol is only secure in the presence of relatively weak semi-honest adversaries. Thus, an important question is how to "convert" the protocol into one that is secure in the presence of malicious adversaries, while preserving the efficiency of the original protocol to the greatest extent possible. Of course, one possibility is to use the compiler of Goldreich, Micali and Wigderson [11]. This compiler converts any protocol that is secure for semi-honest adversaries into one that is secure for malicious adversaries, and as such is a powerful tool for demonstrating feasibility. However, it is based on

reducing the statement that needs to be proved (in our case, the honesty of the parties' behavior) to an NP-complete problem, and using generic zero-knowledge proofs to prove this statement. The resulting secure protocol therefore runs in polynomial time but is rather inefficient. (For more details on existing methods for proving security against malicious behavior see the section on related work below.)

Malicious behavior and cut-and-choose. Consider for a moment what happens if party P_1 is malicious. In such a case, it can construct a garbled circuit that computes a function that is different to the one that P_1 and P_2 agreed to compute. A folklore solution to this problem uses the “cut-and-choose” technique. According to this technique, P_1 first constructs *many* garbled circuits and sends them to P_2 . Then, P_2 asks P_1 to “open” half of them (namely, reveal the decryption keys corresponding to these circuits). P_1 opens the requested half, and P_2 checks that they were constructed correctly. If they were, then P_2 evaluates the rest of the circuits and derives the output from them. The idea behind this methodology is that if a malicious P_1 constructs the circuits incorrectly, then P_2 will detect this with high probability. Clearly, this solution solves the problem of P_1 constructing the circuit incorrectly. However, it does not suffice. First, it creates new problems within itself. Most outstandingly, once the parties now evaluate a number of circuits, some mechanism must be employed to make sure that they use the same input when evaluating each circuit (otherwise, as we show below, an adversarial party could learn more information than allowed). Second, in order to present a proof of security based on *simulation*, there are additional requirements that are not dealt with by just employing cut-and-choose (e.g., input extraction). Third, the folklore description of cut-and-choose is very vague and there are a number of details that are crucial when implementing it. For example, if P_2 evaluates many circuits, then the protocol must specify what P_2 should do if it does not receive the same output in every circuit. If the protocol requires P_2 to abort in this case (because it detected cheating from P_1), then this behavior actually yields a concrete attack in which P_1 can always learn a specified bit of P_2 's input. It can be shown that P_2 must take the majority output and proceed, even if it knows that P_1 has attempted to cheat. This is just one example of a subtlety that must be dealt with. Another example relates to the fact that P_1 may be able to construct a circuit that can be opened with two different sets of keys: the first set opens the circuit correctly and the second incorrectly. In such a case, an adversarial P_1 can pass the basic cut-and-choose test by opening the circuits to be checked correctly. However, it can also supply incorrect keys to the circuits to be computed and thus cause the output of the honest party to be incorrect.

Our contributions. This paper provides several contributions:

- *Efficient protocol for malicious parties:* We present an implementation of Yao's protocol with the cut-and-choose methodology, which is secure in the presence of malicious adversaries and is computationally efficient: the protocol does not use public-key operations, except for performing oblivious

transfers for every input bit of P_2 . For n -bit inputs and a statistical security parameter s the protocol uses $O(\max(s, n))$ oblivious transfers. Thus, when the input is as large as the security parameter, only $O(1)$ oblivious transfers are needed per input bit.

Beyond carefully implementing the cut-and-choose technique on the circuits in order to ensure that the garbled circuits are constructed correctly, we present a new method for enforcing the parties to use the same input in every circuit. This method involves “consistency checks” that are based on cut-and-choose tests which are applied to *sets of commitments* to the garbled values associated with the input wires of the circuit, rather than to the circuits themselves.

In actuality, we combine the cut-and-choose test over the circuits together with the cut-and-choose test over the commitments in order to obtain a secure solution. The test is rather complex conceptually, but is exceedingly simple to implement. Specifically, P_1 just needs to generate a number of commitments to the garbled values associated with the input wires, and then open them based on cut-and-choose queries from P_2 . (Actually, these cut-and-choose queries are chosen jointly by the parties using a simple coin-tossing protocol; this is necessary for achieving simulation.)

We note that the use of cut-and-choose inevitably incurs a higher communication overhead. We also note that in this work we emphasized providing a clear and full proof of the protocol, rather than fully optimizing its overhead at the expense of complicating the proof.

- *Simulation based proof:* We present a *rigorous* proof of the security of the protocol, based on the real/ideal-model simulation paradigm [5,9]. The proof is in the standard model, with no random oracle model or common random string assumptions. The protocol was designed to support such a proof, rather than make do with separate proofs of privacy and correctness. (It is well-known that it is strictly harder to obtain a simulation based proof rather than security under such definitions.) One important advantage of simulation based proofs is that they enable the use of the protocol as a building block in more complicated protocols, while proving the security of the latter using general composition theorems like those of [5,9]. (For example, the secure protocol of [1] for finding the k^{th} ranked element is based on invoking several secure computations of simpler functions, and provides simulation based security against malicious adversaries if the invoked computations have a simulation based proof. However, prior to our work there was no known way, except for the GMW compiler, of efficiently implementing these computations with this level of security.) See [5,9] for more discussion on the importance of simulation-based definitions.
- *A black-box reduction:* Our protocol can be interpreted as a constant-round black-box reduction of secure two-party computation to oblivious transfer and perfectly-hiding commitments. The perfectly-hiding commitments are only used for conducting, in $O(1)$ rounds, joint coin-tossing of a string of

length s , where s is a statistical security parameter. This coin-tossing can be done sequentially (bit by bit), without using perfectly-hiding commitments. We therefore also obtain an $O(s)$ round black-box reduction of secure two-party computation to oblivious transfer alone. These two reductions are comparable to Kilian's reduction, which uses OT alone but incurs a number of rounds which is linear in the depth of the circuit [18]. In addition, our reduction is much more efficient than that of [18].

Related work. As we have mentioned, this paper presents a protocol which **(1)** has a proof of security against malicious adversaries in the standard model, according to the real/ideal model simulation definition, **(2)** has essentially the same computational overhead as Yao's original protocol (which is only secure against semi-honest adversaries), and **(3)** has a somewhat larger communication overhead, which depends on a statistical security parameter s .

We compare this result to other methods for securing Yao's protocol against malicious parties. There are several possible approaches to this task:

- The parties can reduce the statement about the honesty of their behavior to a statement which has a well-known zero-knowledge proof, and then prove this statement. This is the approach taken by the GMW compiler [11]. The resulting secure protocol is not black-box, and is rather inefficient.
- Another approach is to apply a cut-and-choose modification to Yao's protocol. Mohassel and Franklin [23] show such a protocol which has about the same overhead as ours, namely a communication overhead of $O(|C|s + n^2s)$ for a circuit C with n inputs, and a statistical security parameter s . This result was improved by Woodruff [28], who describes how to reduce the communication to $O(|C|s + ns) = O(|C|s)$, using expanders. The protocol of [23] provides output to the circuit evaluator alone. It enables, however, the circuit constructor to carry out the following attack: it can corrupt, say, its OT input which corresponds to a 0 value of the first input bit of the circuit evaluator, while not corrupting the OT input for the 1 value. Other than that it follows the protocol. This behavior forces the circuit evaluator to abort if its first input bit is 0, while if its first input bit is 1 it does not learn anything at all about the attack. If the evaluator complains, then the circuit constructor can conclude that its first input bit is 0, and therefore the evaluator cannot complain if it wants to preserve its privacy. (This attack is similar to the attack we describe in Section 3.2 where we discuss the encoding of P_2 's input.) The protocol therefore does not provide security according to a standard definition. (We note however that this attack can be prevented using the methods we describe in Section 3.2 for encoding P_2 's input.) Another protocol which is based on cut-and-choose is described in [19]. This protocol uses committed OT to address attacks similar to the one described above. We stress that both of these papers ([23,19]) lack a full proof of security, and to our best judgment they need considerable changes in order to support security according to a simulation based definition.
- Jarecki and Shmatikov [15] designed a protocol in which the parties efficiently prove, gate by gate, that their behavior is correct. The protocol is based on

the use of a special homomorphic encryption system, which is used to encode the gates of the table (compared to the use of symmetric encryption in Yao's original protocol and in our paper). The protocol is secure in a universally composable way under the decisional composite residuosity and the strong RSA assumptions, assuming a common reference string.

In this paper, we construct an efficient protocol for *general* secure computation. Thus, we do not (and cannot) compete with protocols that are constructed for specific tasks, like voting, auctions, etcetera. We also do not discuss here the large body of work that considers the efficiency of secure *multi-party* computation.

Organization. We present standard definitions of security for secure two-party computation in Section 2.1. Then, in Section 2.2 we show that a functionality that provides outputs to both parties can be securely reduced to one which provides output for a single party, and therefore we can focus on the latter case. In Section 3 we describe our protocol, prove its security, and analyze its efficiency. The basic protocol we describe increases the number of inputs, and therefore the number of OT invocations. In Section 5.2 we show how to reduce this number of OT invocations in order to improve efficiency. We remark that a description of Yao's basic protocol for two-party computation, secure against semi-honest adversaries, is provided in [20].

2 Preliminaries

2.1 Definitions – Secure Computation

In this section we present the definition for secure two-party computation. The following description and definition is based on [9, Chapter 7], which in turn follows [12,22,4,5].

Two-party computation. A two-party protocol problem is cast by specifying a random process that maps pairs of inputs to pairs of outputs (one for each party). We refer to such a process as a *functionality* and denote it $f : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{0, 1\}^* \times \{0, 1\}^*$, where $f = (f_1, f_2)$. That is, for every pair of inputs (x, y) , the output-pair is a random variable $(f_1(x, y), f_2(x, y))$ ranging over pairs of strings. The first party (with input x) wishes to obtain $f_1(x, y)$ and the second party (with input y) wishes to obtain $f_2(x, y)$.

Adversarial behavior. Loosely speaking, the aim of a secure two-party protocol is to protect an honest party against dishonest behavior by the other party. In this paper, we consider *malicious adversaries* who may arbitrarily deviate from the specified protocol. When considering malicious adversaries, there are certain undesirable actions that cannot be prevented. Specifically, a party may refuse to participate in the protocol, may substitute its local input (and use instead a different input) and may abort the protocol prematurely. One ramification of the

adversary’s ability to abort, is that it is impossible to achieve “fairness”. That is, the adversary may obtain its output while the honest party does not. As is standard for two-party computation, in this work we consider a static corruption model, where one of the parties is adversarial and the other is honest.

Security of protocols (informal). The security of a protocol is analyzed by comparing what an adversary can do in the protocol to what it can do in an ideal scenario that is secure by definition. This is formalized by considering an *ideal* computation involving an incorruptible *trusted third party* to whom the parties send their inputs. The trusted party computes the functionality on the inputs and returns to each party its respective output. Loosely speaking, a protocol is secure if any adversary interacting in the real protocol (where no trusted third party exists) can do no more harm than if it was involved in the above-described ideal computation.

Execution in the ideal model. As we have mentioned, some malicious behavior cannot be prevented (for example, early aborting). This behavior is therefore incorporated into the ideal model. An ideal execution proceeds as follows:

Inputs: Each party obtains an input, denoted w ($w = x$ for P_1 , and $w = y$ for P_2).

Send inputs to trusted party: An honest party always sends w to the trusted party. A malicious party may, depending on w , either abort or send some $w' \in \{0, 1\}^{|w|}$ to the trusted party.

Trusted party answers first party: In case it has obtained an input pair (x, y) , the trusted party first replies to the first party with $f_1(x, y)$. Otherwise (i.e., in case it receives only one valid input), the trusted party replies to both parties with a special symbol \perp .

Trusted party answers second party: In case the first party is malicious it may, depending on its input and the trusted party’s answer, decide to *stop* the trusted party by sending it \perp . In this case the trusted party sends \perp to the second party. Otherwise the trusted party sends $f_2(x, y)$ to the second party.

Outputs: An honest party always outputs the message it has obtained from the trusted party. A malicious party may output an arbitrary (probabilistic polynomial-time computable) function of its initial input and the message obtained from the trusted party.

Let $f : \{0, 1\}^* \times \{0, 1\}^* \mapsto \{0, 1\}^* \times \{0, 1\}^*$ be a functionality, where $f = (f_1, f_2)$, and let $\overline{M} = (M_1, M_2)$ be a pair of non-uniform probabilistic *expected* polynomial-time machines (representing parties in the ideal model). Such a pair is *admissible* if for at least one $i \in \{1, 2\}$ we have that M_i is honest (i.e., follows the honest party instructions in the above-described ideal execution). Then, the joint execution of f under \overline{M} in the ideal model (on input pair (x, y)), denoted $\text{IDEAL}_{f, \overline{M}}(x, y)$, is defined as the output pair of M_1 and M_2 from the above ideal execution.

Execution in the real model. We next consider the real model in which a real (two-party) protocol is executed (and there exists no trusted third party). In this case, a malicious party may follow an arbitrary feasible strategy; that is, any strategy implementable by non-uniform probabilistic polynomial-time machines.

Let f be as above and let Π be a two-party protocol for computing f . Furthermore, let $\overline{M} = (M_1, M_2)$ be a pair of non-uniform probabilistic polynomial-time machines (representing parties in the real model). Such a pair is admissible if for at least one $i \in \{1, 2\}$ we have that M_i is honest (i.e., follows the strategy specified by Π). Then, the joint execution of Π under \overline{M} in the real model (on input pair (x, y)), denoted $\text{REAL}_{\Pi, \overline{M}}(x, y)$, is defined as the output pair of M_1 and M_2 resulting from the protocol interaction.

Security as emulation of a real execution in the ideal model. Having defined the ideal and real models, we can now define security of protocols. Loosely speaking, the definition asserts that a secure two-party protocol emulates the ideal model (in which a trusted party exists). This is formulated by saying that admissible pairs in the ideal model are able to simulate admissible pairs in an execution of a secure real-model protocol.

Definition 1. (secure two-party computation): *Let f and Π be as above. Protocol Π is said to securely compute f (in the malicious model) if for every pair of admissible non-uniform probabilistic polynomial-time machines $\overline{A} = (A_1, A_2)$ for the real model, there exists a pair of admissible non-uniform probabilistic expected polynomial-time machines $\overline{B} = (B_1, B_2)$ for the ideal model, such that*

$$\left\{ \text{IDEAL}_{f, \overline{B}}(x, y) \right\}_{x, y \text{ s.t. } |x|=|y|} \stackrel{c}{\equiv} \left\{ \text{REAL}_{\Pi, \overline{A}}(x, y) \right\}_{x, y \text{ s.t. } |x|=|y|}$$

Namely, the two distributions are computationally indistinguishable.

We note that the above definition assumes that the parties know the input lengths (this can be seen from the requirement that $|x| = |y|$). Some restriction on the input lengths is unavoidable, see [9, Section 7.1] for discussion. We also note that we allow the ideal adversary/simulator to run in expected (rather than strict) polynomial-time. This is essential for achieving constant-round protocols; see [3].

We denote the security parameter by n and, for the sake of simplicity, unify it with the length of the inputs (thus we consider security for “all sufficiently long inputs”). Everything in the paper remains the same if a separate security parameter n is used, and we consider security for inputs of all lengths. We will also use a statistical security parameter s ; see the beginning of Section 3.1 for an explanation of the use of this separate parameter.

The hybrid model. Our protocol uses a secure oblivious transfer protocol as a subprotocol. It has been shown in [5] that it suffices to analyze the security of such a protocol in a hybrid model in which the parties interact with each other *and*

have access to a trusted party that computes the oblivious transfer protocol for them. We remark that the composition theorem of [5] holds for the case that the subprotocol executions are all run sequentially (and the messages of the protocol calling the subprotocol do not overlap with any execution). We also remark that if the oblivious transfer subprotocol is secure under parallel composition, then it is straightforward to extend [5] so that the subprotocols may be run in parallel (again, as long as the messages of the protocol calling the subprotocol do not overlap with any execution).

2.2 Functionalities That Provide Output to a Single Party

In the definition above, we have considered the case that both parties receive output, and these outputs may be *different*. However, the presentation of our protocol is far simpler for the case that only party P_2 receives output. We will show now that this suffices for the general case. That is, any protocol that can securely compute *any* efficient functionality $f(x, y)$ where only P_2 receives output, can be used to securely compute *any* efficient functionality $f = (f_1, f_2)$ where party P_1 receives $f_1(x, y)$ and party P_2 receives $f_2(x, y)$.

Let $f = (f_1, f_2)$ be a functionality. We wish to construct a secure protocol in which P_1 receives $f_1(x, y)$ and P_2 receives $f_2(x, y)$; as a building block we use a protocol for computing any efficient functionality with the limitation that only P_2 receives output. Let \mathcal{F} be a field that contains the range of values $\{f_1(x, y)\}_{x, y \in \{0, 1\}^n}$, and let p, a, b be randomly chosen elements in \mathcal{F} . Then, in addition to x , party P_1 's input includes the elements p, a, b . Furthermore, define a functionality g (that has only a single output) as follows: $g((p, a, b, x), y) = (\alpha, \beta, f_2(x, y))$, where $\alpha = p + f_1(x, y)$, $\beta = a \cdot \alpha + b$, and the arithmetic operations are defined in \mathcal{F} . Note that α is a one-time pad encryption of P_1 's output $f_1(x, y)$, and β is an information-theoretic message authentication tag of α (specifically, $a\alpha + b$ is a pairwise-independent hash of α). Now, the parties compute the functionality g , using a secure protocol in which only P_2 receives output. Following this, P_2 sends the pair (α, β) to P_1 . Party P_1 checks that $\beta = a \cdot \alpha + b$; if yes, it outputs $\alpha - p$, and otherwise it outputs \perp .

It is easy to see that P_2 learns nothing about P_1 's output $f_1(x, y)$, and that it cannot alter the output that P_1 will receive (beyond causing it to abort), except with probability $1/|\mathcal{F}|$. (We assume that $1/|\mathcal{F}|$ is the required probability for detecting attempts to alter the output. If it is required instead that any change by P_2 to P_1 's output is detected with probability 2^{-s} , then the parameters a, b and the computation of $\beta = a \cdot \alpha + b$ can be defined in a field whose representation is s bits long.) We remark that it is also straightforward to construct a simulator for the above protocol. (Note that in order to meet Definition 1, one must actually switch the roles of P_1 and P_2 above.)

We remark that the circuit for computing g is only mildly larger than that for computing f . Thus, the construction above is also efficient and has only a mild effect on the complexity of the secure protocol.

3 The Protocol

Our protocol is based upon Yao’s garbled circuit construction, which is secure in the presence of semi-honest adversaries [29]. That protocol has two parties: P_1 (who is the the *sender*, or circuit *constructor*), and P_2 (who is the *receiver*, or the circuit *evaluator*). The protocol is described and proved in [20]. Our presentation from here on assumes full familiarity with Yao’s basic protocol.

There are a number of issues that must be dealt with when attempting to make Yao’s protocol secure against malicious adversaries rather than just semi-honest ones (beyond the trivial observation that the oblivious transfer subprotocol must now be secure in the presence of malicious adversaries).

First and foremost, a malicious P_1 must be forced to construct the garbled circuit correctly so that it indeed computes the desired function. The method that is typically referred to for this task is called *cut-and-choose*. According to this methodology, P_1 constructs many independent copies of the garbled circuit and sends them to P_2 . Party P_2 then asks P_1 to open half of them (chosen randomly). After P_1 does so, and party P_2 checks that the opened circuits are correct, P_2 is convinced that most of the remaining (unopened) garbled circuits are also constructed correctly. (If there are many incorrectly constructed circuits, then with high probability, one of those circuits will be in the set that P_2 asks to open.) The parties can then evaluate the remaining unopened garbled circuits as in the original protocol for semi-honest adversaries, and take the majority output-value.¹

The cut-and-choose technique described above indeed solves the problem of a malicious P_1 constructing incorrect circuits. However, it also generates new problems! The primary problem that arises is that since there are now many circuits being evaluated, we must make sure that both P_1 and P_2 use the same inputs in each circuit; we call these *consistency checks*. (Consistency checks are important since if the parties were able to provide different inputs to different copies of the circuit, then they can learn information that is different from the desired output of the function. It is obvious that P_2 can do so, since it observes the outputs of all circuits, but in fact even P_1 , who only gets to see the majority

¹ The reason for taking the majority value as the output is that the aforementioned test only reveals a single incorrectly constructed circuit with probability $1/2$. Therefore, if P_1 generates a single or constant number of “bad” circuits, there is a reasonable chance that it will not be caught. In contrast, there is only an exponentially small probability that the test reveals no corrupt circuit *and* at the same time a majority of the circuits that are not checked are incorrect. Consequently, with overwhelming probability it holds that if the test succeeds and P_2 takes the majority result of the remaining circuits, the result is correct. We remark that the alternative of aborting in case not all the outputs are the same (namely, where cheating is detected) is not secure and actually yields a concrete attack. The attack works as follows. Assume that P_1 is corrupted and that it constructs all of the circuits correctly except for one. The “incorrect circuit” is constructed so that it computes the exclusive-or of the desired function f with the first bit of P_2 ’s input. Now, if P_2 ’s policy is to abort as soon as two outputs are not the same then P_1 learns the first bit of P_2 ’s input.

output, can learn additional information: information².) Another problem that arises when proving security is that the simulator must be able to fool P_2 and give it incorrect circuits (even though P_2 runs a cut-and-choose test). This is solved using rather standard techniques, like choosing the circuits to be opened via a coin-tossing protocol (to our knowledge, this issue has gone unnoticed in all previous applications of cut-and-choose to Yao's protocol). Yet another problem is that P_1 might provide corrupt inputs to some of P_2 's possible choices in the OT protocols. P_1 might then learn P_2 's input based on whether or not P_2 aborts the protocol.

We begin by presenting a high-level overview of the protocol. We then proceed to describe the consistency checks, and finally the full protocol.

3.1 High-Level Overview

We work with two security parameters. The parameter n is the security parameter for the commitment schemes, encryption, and the oblivious transfer protocol. The parameter s is a statistical security parameter which specifies how many garbled circuits are used. The difference between these parameters is due to the fact that the value of n depends on computational assumptions, whereas the value of s reflects the possible error probability that is incurred by the cut-and-choose technique and as such is a "statistical" security parameter. Although it is possible to use a single parameter n , it may be possible to take s to be much smaller than n . Recall that for simplicity, and in order to reduce the number of parameters, we denote the length of the input by n as well.

Protocol 1. (high-level overview): *Parties P_1 and P_2 have respective inputs x and y , and wish to compute the output $f(x, y)$ for P_2 .*

0. *The parties decide on a circuit computing f . They then change the circuit by replacing each input wire of P_2 by a gate whose input consists of s new input wires of P_2 and whose output is the exclusive-or of these wires (such an s -bit exclusive-or gate can be implemented using $s-1$ two-bit exclusive-or gates). Consequently, the number of input wires of P_2 increases by a factor of s . (In Section 5.2, we show how to reduce the number of inputs.)*
1. *P_1 commits to s different garbled circuits computing f , where s is a statistical security parameter. P_1 also generates additional commitments to the garbled values corresponding to the input wires of the circuits. These commitments are constructed in a special way in order to enable consistency checks.*
2. *For every input bit of P_2 , parties P_1 and P_2 run a 1-out-of-2 oblivious transfer protocol in which P_2 learns the garbled values of input wires corresponding to its input.*

² Suppose, for example, that the protocol computes n invocations of a circuit computing the inner-product between n bit inputs. A malicious P_2 could provide the inputs $\langle 10 \cdots 0 \rangle, \langle 010 \cdots 0 \rangle, \dots, \langle 0 \cdots 01 \rangle$, and learn all of P_1 's input. If, on the other hand, P_1 is malicious, it could also provide the inputs $\langle 10 \cdots 0 \rangle, \langle 010 \cdots 0 \rangle, \dots, \langle 0 \cdots 01 \rangle$. In this case, P_2 sends it the value which is output by the majority of the circuits, and which is equal to the majority value of P_2 's input bits.

3. P_1 sends to P_2 all the commitments of Step 1.
4. P_1 and P_2 run a coin-tossing protocol in order to choose a random string that defines which commitments and which garbled circuits will be opened.
5. P_1 opens the garbled circuits and committed input values that were chosen in the previous step. P_2 verifies the correctness of the opened circuits and runs consistency checks based on the decommitted input values.
6. P_1 sends P_2 the garbled values corresponding to P_1 's input wires in the unopened circuits. P_2 runs consistency checks on these values as well.
7. Assuming that all of the checks pass, P_2 evaluates the unopened circuits and takes the majority value as its output.

3.2 Checks for Correctness and Consistency

As can be seen from the above overview, P_1 and P_2 run a number of checks, with the aim of forcing a potentially malicious P_1 to construct the circuits correctly and use the same inputs in (most of) the evaluated circuits. This section describes these checks.

Encoding P_2 's input: As mentioned above, a malicious P_1 may provide corrupt input to one of P_2 's possible inputs in an OT protocol. If P_2 chooses to learn this input it will not be able to decode the garbled tables which use this value, and it will therefore have to abort. If P_2 chooses to learn the other input associated with this wire then it will not notice that the first input is corrupt. P_1 can therefore learn P_2 's input based on whether or not P_2 aborts. (Note that checking that the circuit is well-formed will not help in thwarting this attack, since the attack is based on changing P_1 's input to the OT protocol.) The attack is prevented by the parties replacing each input bit of P_2 with s new input bits whose exclusive-or is used instead of the original input (this step was described as Step 0 of Protocol 1). P_2 therefore has 2^{s-1} ways to encode a 0 input, and 2^{s-1} ways to encode a 1, and given its input it chooses the encoding to use with uniform probability. The parties then execute the protocol with the new circuit, and P_2 uses oblivious transfer to learn the garbled values of its new inputs. As is shown in the full paper, if P_1 supplies incorrect values as garbled values that are associated with P_2 's input, the probability of P_2 detecting this cheating is *almost independent* (up to a bias of 2^{-s+1}) of P_2 's actual input. This is not true if P_2 's inputs are not "split" in the way described above. The encoding presented here increases the number of P_2 's input bits and, respectively, the number of OTs, from n to ns . In Section 5.2 we show how to reduce the number of new inputs for P_2 (and thus OTs) to a total of only $O(\max(s, n))$.

An unsatisfactory method for proving consistency of P_1 's input: Consider the following idea for forcing P_1 to provide the same input to all circuits. Let s be a security parameter and assume that there are s garbled copies of the circuit. Then, P_1 generates two ordered sets of commitments for every wire of the circuit. Each set contains s commitments: the "0 set" contains commitments to the garbled encodings of 0 for this wire in every circuit, and the "1 set" contains

commitments to the garbled encodings of 1 for this wire in every circuit. P_2 receives these commitments from P_1 and then chooses a random subset of the circuits, which will be defined as *check-circuits*. These circuits will never be evaluated and are used only for checking correctness and consistency. Specifically, P_2 asks P_1 to de-garble all of the check-circuits and to open the values that correspond to the check-circuits in *both* commitment sets. (That is, if circuit i is a check-circuit, then P_1 decommits to both the 0 encoding and 1 encoding of all the input wires in circuit i .) Upon receiving the decommitments, P_2 verifies that all opened commitments from the “0 set” correspond to garbled values of 0, and that a similar property holds for commitments from the “1 set”.

It now remains for P_2 to evaluate the remaining circuits. In order to do this, P_1 provides (for each of its input wires) the garbled values that are associated with the wire in all of the remaining circuits. Then, P_1 must prove that all of these values come from the same set, without revealing whether the set that they come from is the “0 set” or the “1 set” (otherwise, P_2 will know P_1 ’s input). In this way, on the one hand, P_2 does not learn the input of P_1 , and on the other hand, it is guaranteed that all of the values come from the same set, and so P_1 is forced into using the same input in all circuits. This proof can be carried out using, for example, the proofs of partial knowledge of [6]. However, this would require n proofs, each for s values, thereby incurring $O(ns)$ costly asymmetric operations which we want to avoid.

Proving consistency of P_1 ’s input: P_1 can prove consistency of its inputs without using public-key operations. The proof is based on a cut-and-choose test for the consistency of the commitment sets, which is combined with the cut-and-choose test for the correctness of the circuits. (Note that in the previous proposal, there is only one cut-and-choose test, and it is for the correctness of the circuits.) We start by providing a high level description of the proof of consistency: The proof is based on P_1 constructing, *for each of its input wires*, s pairs of sets of commitments. One set in every pair contains commitments to the 0 values of this wire in *all* circuits, and the other set is the same with respect to 1. The protocol chooses a random subset of these pairs, and a random subset of the circuits, and checks that these sets provide consistent inputs for these circuits. Then the protocol evaluates the *remaining* circuits, and asks P_1 to open, in each of the *remaining* pairs, and only in one set in every pair, its garbled values for all evaluated circuits. (In this way, P_2 does not learn whether these garbled values correspond to a 0 or to a 1.) In order for the committed sets and circuits to pass P_2 ’s checks, there must be large subsets C and S , of the circuits and commitment sets, respectively, such that every choice of a circuit from C and a commitment set from S results in a circuit and garbled values which compute the desired function f . P_2 accepts the verification stage only if all the circuits and sets it chooses to check are from C and S , respectively. This means that if P_2 does not abort then circuits which are not from C are likely to be a minority of the evaluated circuits, and a similar argument holds for S . Therefore the majority result of the evaluation stage is correct. The exact construction is as follows:

STAGE 1 – COMMITMENTS: P_1 generates s garbled versions of the circuit. Furthermore, it generates commitments to the garbled values of the wires corresponding to P_2 's input in each circuit. These commitments are generated in ordered pairs so that the first item in a pair corresponds to the 0 value and the second to the 1 value. The procedure regarding the input bits of P_1 is more complicated (see Figure 1 for a diagram explaining this construction). P_1 generates s pairs of sets of committed values for each of its input wires. Specifically, for every input wire i of P_1 , it generates s sets of the form $\{W_{i,j}, W'_{i,j}\}_{j=1}^s$; we call these commitment sets. Before describing the content of these sets, denote by $k_{i,r}^b$ the garbled value that is assigned to the value $b \in \{0,1\}$ in wire i of circuit r . Then, the sets $W_{i,j}$ and $W'_{i,j}$ both contain $s + 1$ commitments and are defined as follows. Let $b \in_R \{0,1\}$ be a random bit, chosen independently for every $\{W_{i,j}, W'_{i,j}\}$ pair. Define $W_{i,j}$ to contain a commitment to b , as well as commitments to the garbled value corresponding to b in wire i in all of the s circuits, and define $W'_{i,j}$ similarly, but with respect to $1 - b$. In other words, $W_{i,j} = \{\text{com}(b), \text{com}(k_{i,1}^b), \dots, \text{com}(k_{i,s}^b)\}$ and $W'_{i,j} = \{\text{com}(1-b), \text{com}(k_{i,1}^{1-b}), \dots, \text{com}(k_{i,s}^{1-b})\}$. We stress that in each of the pairs $(W_{i,1}, W'_{i,1}), \dots, (W_{i,s}, W'_{i,s})$, the values that are committed to are the same. The only difference is that independent randomness is used in each pair for choosing b and constructing the commitments. We call the first bit committed to in a commitment set the indicator bit.

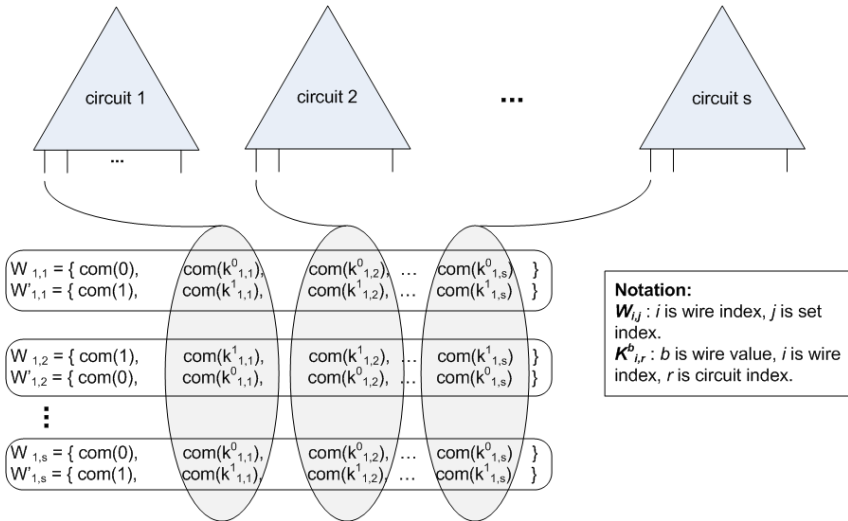


Fig. 1. The commitment sets corresponding to P_1 's first input wire

After constructing these circuits and commitment sets, party P_1 sends to P_2 all of the s garbled circuits (i.e., the garbled gate-tables and output-tables, but *not* the garbled values corresponding to the input wires), and all the commitment

sets. Note that if P_1 's input is of length n , then there are sn pairs of commitment sets; and a total of $sn(2s + 2) = O(s^2n)$ commitments.

STAGE 2 – CHALLENGE: Two random strings $\rho, \rho' \in_R \{0, 1\}^s$ are chosen and sent to P_1 (in the actual protocol, these strings are determined via a simple coin-tossing protocol). The string ρ is a challenge indicating which garbled circuits to open, and the string ρ' is a challenge indicating which commitment sets to open. We call the opened circuits check-circuits and the unopened ones evaluation-circuits. Likewise, we call the opened sets check-sets and the unopened ones evaluation-sets. A circuit (resp., commitment set) is defined to be a check-circuit (resp., check-set) if the corresponding bit in ρ (resp., ρ') equals 1; otherwise, it is defined to be an evaluation-circuit (resp., evaluation-set).

STAGE 3 – OPENING: First, party P_1 opens all the commitments corresponding to P_2 's input wires in all of the check-circuits. Second, in all of the *check-sets* P_1 opens the commitments that correspond to *check-circuits*. That is, if circuit r is a check circuit, then P_1 decommits to all of the values $\text{com}(k_{i,r}^0), \text{com}(k_{i,r}^1)$ in check-sets, where i is any of P_1 's input bits. Finally, for every check-set, P_1 opens the commitment to the indicator bit, the initial value in each of the sets $W_{i,j}, W'_{i,j}$. See Figure 2 for a diagram in which the values which are opened are highlighted (the diagram refers to only one of P_1 's input wires in the circuit).

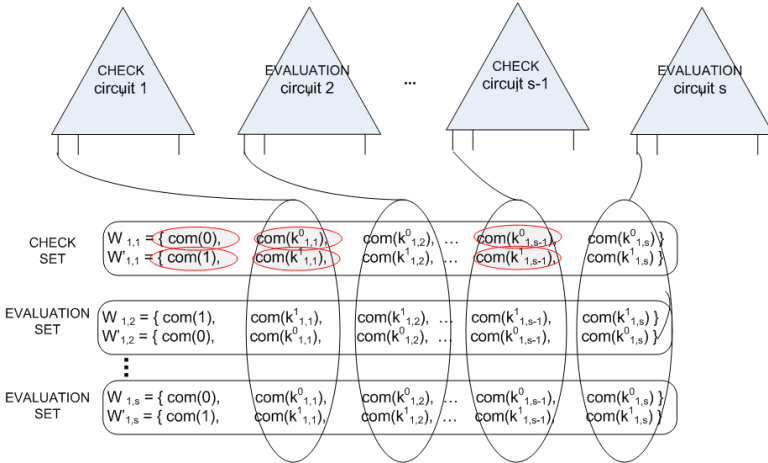


Fig. 2. In every check-set, the commitment to the indicator bit, and the commitments corresponding to check-circuits are all opened.

STAGE 4 – VERIFICATION: In this step, party P_2 verifies that all of the check-circuits were correctly constructed. In addition, it verifies that all of the opened commitments in sets whose first item is a commitment to 0 are to garbled encodings of 0; likewise for 1. These checks are carried out as follows. First, in all of the check-circuits, P_2 receives the decommitments to the garbled values

corresponding to its own input, and by the order of the commitments P_2 knows which value corresponds to 0 and which value corresponds to 1. Second, for every check-circuit, P_2 receives decommitments to the garbled input values of P_1 in all the check-sets, along with a bit indicating whether these garbled values correspond to 0 or to 1. It first checks that for every wire, the garbled values of 0 (resp., of 1) are all equal. Then, the above decommitments enable the complete opening of the garbled circuits (i.e., the decryption of all of the garbled tables). Once this has been carried out, it is possible to simply check that the check-circuits are all correctly constructed. Namely, that they agree with a specific and agreed-upon circuit computing f .

STAGE 5 – EVALUATION AND VERIFICATION: Party P_1 reveals the garbled values corresponding to its input: If i is a wire that corresponds to a bit of P_1 's input and r is an evaluation-circuit, then P_1 decommits to the commitments $k_{i,r}^b$ in all of the *evaluation-sets*, where b is the value of its input bit. This is depicted in Figure 3. Finally, P_2 verifies that (1) for every input wire, all of the opened commitments that were opened in evaluation-sets contain the same garbled value, and (2) for every i, j P_1 opened commitments of evaluated circuits in exactly one of $W_{i,j}$ or $W'_{i,j}$. If these checks pass, it continues to evaluate the circuit.

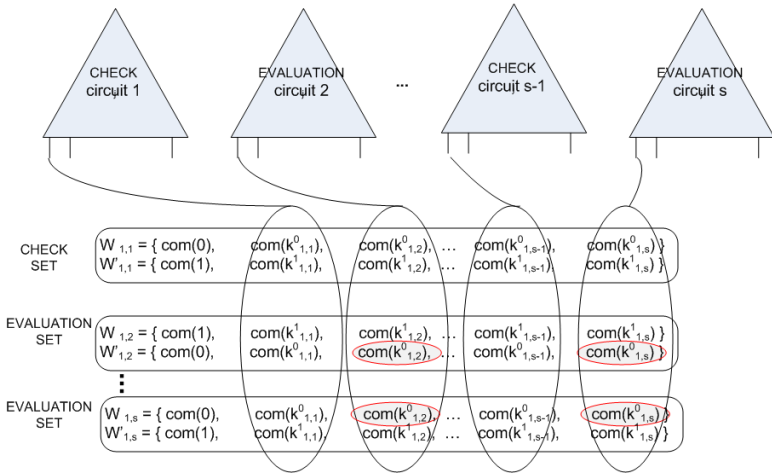


Fig. 3. P_1 opens in the evaluation-sets, the commitments that correspond to its input. In every evaluation-set these commitments come from the same item in the pair.

Intuition. Having described the mechanism for checking consistency, we now provide some intuition as to why it is correct. A simple cut-and-choose check verifies that most of the evaluated circuits are correctly constructed. The main remaining issue is ensuring that P_1 's inputs to most circuits are consistent. If P_1 wants to provide different inputs to a certain wire in two circuits, then *all* the $W_{i,j}$ (or $W'_{i,j}$) sets it opens in *evaluation-sets* must contain a commitment to 0 in the first circuit and a commitment to 1 in the other circuit. However,

if any of these sets is chosen to be checked, and the circuits are among the checked circuits, then P_2 aborts. This means that if P_1 attempts to provide different inputs to two circuits and they are checked, it is almost surely caught. Now, since P_2 outputs the majority output of the evaluated circuits, the result is affected by P_1 providing different inputs only if these inputs affect a constant fraction of the circuits. But since all of these circuits must not be checked, P_1 's probability of success is exponentially small in s .

3.3 The Full Protocol

We now describe the full protocol in detail. We use the notation com_b to refer to a perfectly binding commitment scheme, and com_h to refer to a perfectly hiding commitment scheme (See [8] for definitions).

Protocol 2. (protocol for computing $f(x, y)$):

- **Input:** P_1 has input $x \in \{0, 1\}^n$ and P_2 has input $y \in \{0, 1\}^n$.
- **Auxiliary input:** a statistical security parameter s and the description of a circuit C^0 such that $C^0(x, y) = f(x, y)$.
- **Specified output:** party P_2 should receive $f(x, y)$ and party P_1 should receive no output. (Recall that this suffices for the general case where both parties receive possibly different outputs; see Section 2.2.)
- **The protocol**
 0. **CIRCUIT CONSTRUCTION:** The parties replace C^0 with a circuit C which is constructed by replacing each input wire of P_2 by the result of an exclusive-or of s new input wires of P_2 . (We show in Section 5.2 how the number of new input bits can be reduced.) The number of input wires of P_2 is increased from $|y| = n$ to sn . Let the bit-wise representation of P_2 's original input be $y = y_1 \dots y_n$. Denote its new input as $\hat{y} = \hat{y}_1, \dots, \hat{y}_{ns}$. P_2 chooses its new input at random subject to the constraint $y_i = \hat{y}_{(i-1) \cdot s + 1} \oplus \dots \oplus \hat{y}_{i \cdot s}$.
 1. **COMMITMENT CONSTRUCTION:** P_1 constructs the circuits and commits to them, as follows:
 - (a) P_1 constructs s independent copies of a garbled circuit of C , denoted GC_1, \dots, GC_s .
 - (b) P_1 commits to the garbled values of the wires corresponding to P_2 's input to each circuit. That is, for every input wire i corresponding to an input bit of P_2 , and for every circuit GC_r , P_1 computes the ordered pair $(\text{com}_b(k_{i,r}^0), \text{com}_b(k_{i,r}^1))$, where $k_{i,r}^b$ is the garbled value associated with b on input wire i in circuit GC_r .
 - (c) P_1 computes commitment-sets for the garbled values that correspond to its own inputs to the circuits. That is, for every wire i that corresponds to an input bit of P_1 , it generates s pairs of commitment sets $\{W_{i,j}, W'_{i,j}\}_{j=1}^s$, in the following way:

Denote by $k_{i,r}^b$ the garbled value that was assigned by P_1 to the value $b \in \{0, 1\}$ of wire i in GC_r . Then, P_1 chooses $b \in_R \{0, 1\}$ and computes

$$W_{i,j} = \langle \text{com}_b(b), \text{com}_b(k_{i,1}^b), \dots, \text{com}_b(k_{i,s}^b) \rangle, \quad \text{and}$$

$$W'_{i,j} = \langle \text{com}_b(1-b), \text{com}_b(k_{i,1}^{1-b}), \dots, \text{com}_b(k_{i,s}^{1-b}) \rangle$$

For each i, j , the sets are constructed using independent randomness, and in particular the value of b is chosen independently for every $j = 1 \dots s$. There are a total of ns commitment-sets. We divide them into s supersets, where superset S_j is defined as $S_j = \{(W_{1,j}, W'_{1,j}), \dots, (W_{n,j}, W'_{n,j})\}$. Namely, S_j is the set containing the j^{th} commitment set for all wires.

2. OBLIVIOUS TRANSFERS: For every input bit of P_2 , parties P_1 and P_2 run a 1-out-of-2 oblivious transfer protocol in which P_2 receives the garbled values for the wires that correspond to its input bit (in every circuit). That is, let $c_{i,r}^b$ denote the commitment to the garbled value $k_{i,r}^b$ and let $dc_{i,r}^b$ denote the decommitment value for $c_{i,r}^b$. Furthermore, let i_1, \dots, i_{ns} be the input wires that correspond to P_2 's input.

Then, for every $j = 1, \dots, ns$, parties P_1 and P_2 run a 1-out-of-2 OT protocol in which:

- (a) P_1 's input is the pair of vectors $([dc_{i_j,1}^0, \dots, dc_{i_j,s}^0], [dc_{i_j,1}^1, \dots, dc_{i_j,s}^1])$.
- (b) P_2 's input is its j^{th} input bit \hat{y}_j (and its output should thus be $[dc_{i_j,1}^{\hat{y}_j}, \dots, dc_{i_j,s}^{\hat{y}_j}]$).

If the oblivious transfer protocol provides security for parallel execution, then these executions are run in parallel. Otherwise, they are run sequentially.

3. SEND CIRCUITS AND COMMITMENTS: P_1 sends to P_2 the garbled circuits (i.e., the gate and output tables), as well as all of the commitments that it prepared above.
4. PREPARE CHALLENGE STRINGS: (1) P_2 chooses a random string $\rho_2 \in_R \{0, 1\}^s$ and sends $\text{com}_h(\rho_2)$ to P_1 . (2) P_1 chooses a random string $\rho_1 \in \{0, 1\}^s$ and sends $\text{com}_b(\rho_1)$ to P_2 . (3) P_2 decommits, revealing ρ_2 . (4) P_1 decommits, revealing ρ_1 . (5) P_1 and P_2 set $\rho = \rho_1 \oplus \rho_2$. The above steps are run a second time, defining an additional string ρ' .³

5. DECOMMITMENT PHASE FOR CHECK-CIRCUITS: From here on, we refer to the circuits for which the corresponding bit in ρ is 1 as check-circuits,

³ Recall that ρ and ρ' are used to ensure that P_1 constructs the circuits correctly and uses consistent input in each circuit. Thus, it may seem strange that they are generated via a coin-tossing protocol, and not just chosen singlehandedly by P_2 . Indeed, in order to prove the security of the protocol when P_1 is corrupted, there is no need for a coin-tossing protocol here. However, having P_2 choose ρ and ρ' singlehandedly creates a problem for the simulation in the case that P_2 is corrupted. We therefore use a coin-tossing protocol instead.

and we refer to the other circuits as evaluation-circuits. Likewise, if the j^{th} bit of ρ' equals 1, then the commitments sets in $S_j = \{(W_{i,j}, W'_{i,j})\}_{i=1\dots n}$ are referred to as check-sets; otherwise, they are referred to as evaluation-sets.

For every check-circuit GC_r , party P_1 operates in the following way:

- (a) For every input wire i corresponding to an input bit of P_2 , party P_1 decommits to the pair $(\text{com}(k_{i,r}^0), \text{com}(k_{i,r}^1))$ (namely to both of P_2 's inputs).
- (b) For every input wire i corresponding to an input bit of P_1 , party P_1 decommits to the appropriate values in the superset S_j , in the check-sets $\{W_{i,j}, W'_{i,j}\}$. Specifically, P_1 decommits to the $\text{com}(k_{i,r}^0)$ and $\text{com}(k_{i,r}^1)$ values in $(W_{i,j}, W'_{i,j})$, for every check-set S_j (see Figure 2).

For every pair of check-sets $(W_{i,j}, W'_{i,j})$, party P_1 decommits to the first value in each set (i.e., to the value that is supposed to be a commitment to the indicator bit, $\text{com}(0)$ or $\text{com}(1)$).

6. DECOMMITMENT PHASE FOR P_1 'S INPUT IN EVALUATION-CIRCUITS: P_1 decommits to the garbled values that correspond to its inputs in evaluation-circuits. Let i be the index of an input wire that corresponds to P_1 's input (the following procedure is applied to all such wires). Let b be the binary value that P_1 assigns to input wire i . In every evaluation-set $(W_{i,j}, W'_{i,j})$, P_1 chooses the set (out of $(W_{i,j}, W'_{i,j})$), which corresponds to the value b . It then opens in this set the commitments that correspond to evaluation-circuits. Namely, to the values $k_{i,r}^b$, where r is an index of an evaluation circuit (see Figure 3).
7. CORRECTNESS AND CONSISTENCY CHECKS: P_2 performs the following checks; if any of them fails it aborts.
 - (a) Checking correctness of the check-circuits: P_2 verifies that each check-circuit GC_i is a garbled version of C . This check is carried out by P_2 first constructing the input tables that associate every garbled value of an input wire to a binary value. The input tables for P_2 's inputs are constructed by checking that the decommitments to the pairs $(\text{com}(k_{i,r}^0), \text{com}(k_{i,r}^1))$ (where i is a wire index and r is a circuit index) are valid, and then interpreting the first value to be associated with 0 and the second value to be associated with 1.

Next, P_2 checks the decommitments to P_1 's inputs. This check involves first checking that the decommitment values are valid. Then, P_2 verifies that in each pair of check-sets, one of $(W_{i,j}, W'_{i,j})$ begins with a commitment to 0 (henceforth the 0-tuple), and the other begins with a commitment to 1 (henceforth the 1-tuple). Then P_2 checks that for every wire, the values that are decommitted to in the 0-tuples in all check-sets are all equal, and that a similar property holds for the 1-tuples. P_2 then assigns the logical value of 0 to all of the opened

commitments in the 0-tuples, and the logical value of 1 to the opened commitments in the 1-tuples.

Finally, given all the garbled values to the input wires and their associated binary values, P_2 decrypts the circuit and compares it with the circuit C .

- (b) Verifying P_2 's input in the check-circuits: P_2 verifies that P_1 's decommitments to the wires corresponding to P_2 's input values in the check-circuits are correct, and agree with the logical values of these wires (the indicator bits). P_2 also checks that the inputs it learned in the oblivious transfer stage for the check-circuits correspond to its actual input. Specifically, it checks that the decommitment values that it received in the oblivious transfer stage open the committed values that correspond to the garbled values of its logical input (namely, that it received the first value in the pair if the input bit is 0 and the second value if it is 1).⁴
- (c) Checking P_1 's input to evaluation-circuits: Finally, P_2 verifies that for every input wire i of P_1 the following two properties hold:
- i. In every evaluation-set P_1 chose one of the two sets and decommitted to all the commitments in it which correspond to evaluation-circuits.
 - ii. For every evaluation-circuit, all of the commitments that P_1 opened in evaluation-sets are for the same garbled value.
8. CIRCUIT EVALUATION: If any of the above checks fails, P_2 aborts and outputs \perp . Otherwise, P_2 evaluates the evaluation circuits (in the same way as for the semi-honest protocol of Yao). It might be that in certain circuits the garbled values provided for P_1 's inputs, or the garbled values learned by P_2 in the OT stage, do not match the tables and so decryption of the circuit fails. In this case P_2 also aborts and outputs \perp . Otherwise, P_2 takes the output that appears in most circuits, and outputs it (the proof shows that this value is well defined).

4 Proof of Security

The security of Protocol 2 is stated in the following theorem.

Theorem 2. *Let $f : \{0, 1\}^* \times \{0, 1\}^* \rightarrow \{0, 1\}^*$ be any probabilistic polynomial-time two-party functionality and consider the instantiation of Protocol 2 for functionality f . Assume that the oblivious transfer protocol is secure, that com_b is a perfectly-binding commitment scheme, that com_h is a perfectly-hiding commitment scheme, and that the garbled circuits are constructed as in [20]. Then, Protocol 2 securely computes f .*

⁴ This check is crucial and thus the order of first running the oblivious transfer and then sending the circuits and commitments is not at all arbitrary.

The theorem is proved in two stages: first for the case that P_1 is corrupted and next for the case that P_2 is corrupted. The proof is provided in the full version of this paper. We highlight here the basic intuition behind the proof.

Security against a Malicious P_1 . The proof constructs an ideal-model adversary/simulator which has access to P_1 and to the trusted party, and can simulate the view of an actual run of the protocol. It uses the fact that the strings ρ, ρ' , which choose the circuits and commitment sets that are checked, are uniformly distributed even if P_1 is malicious. The simulator runs the protocol until P_1 opens the commitments of the checked circuits and checked commitment sets, and then rewinds the execution and runs it again with new random ρ, ρ' values. We expect that about one quarter of the circuits are checked in the first execution *and* evaluated in the second execution. For these circuits, in the first execution the simulator learns the translation between the garbled values of P_1 's input wires and the actual values of these wires, and in the second execution it learns the garbled values that are associated with P_1 's input (this association is learned from the garbled values that P_1 sends to P_2). Combining the two, it learns P_1 's input x , which can then be sent to the trusted party. The trusted party answers with $f(x, y)$, which we use to define P_2 's output and complete the simulation.

Security against a Malicious P_2 . Intuitively, the security in this case is derived from the fact that: **(a)** the oblivious transfer protocol is secure, and so P_2 only learns a single set of keys (corresponding to a single input y) for decrypting the garbled circuits, and **(b)** the commitment schemes are hiding and so P_2 does not know what input corresponds to the garbled values that P_1 sends it for evaluating the circuit. Of course, in order to formally prove security we construct an ideal-model simulator B_2 working with an adversary A_2 that has corrupted P_2 . The simulator first extracts A_2 's input bits from the oblivious transfer protocol, and then sends the input y it obtained to the trusted party and receives back $z = f(x, y)$. Given the output, the simulator constructs the garbled circuits. However, rather than constructing them all correctly, for each circuit it tosses a coin and, based on the result, either constructs the circuit correctly, or constructs it to compute the constant function outputting z (the output is received from the trusted party). In order to make sure that the simulator is not caught cheating, it biases the coin-tossing phase so that all of the correctly-constructed garbled circuits are check-circuits, and all of the other circuits are evaluation-circuits (this is why the protocol uses joint coin-tossing rather than let P_2 alone choose the circuits to be opened). A_2 then checks the correctly-constructed circuits, and is satisfied with the result as if it were interacting with a legitimate P_1 . A_2 therefore continues the execution with the circuits which always output z .

5 Efficiency of the Protocol

We discuss below the efficient implementation of the different building blocks of the protocol (namely, encryption, commitment schemes, and oblivious transfer).

The overhead of the protocol depends on a statistical security parameter s . The security proof shows that the adversary's cheating probability is exponentially small in s . We note that in this paper we preferred to present a full and clear proof, rather than overly optimize the construction at the cost of complicating the proof. We have not analyzed the exact constants affecting the dependence of the error probability on the security parameter s .

The communication overhead of the protocol is dominated by sending s copies of the garbled circuit, and $2s(s+1)$ commitments for each of the n inputs of P_1 . In the protocol, the original circuit C^0 is modified by replacing each of the n original input bits of P_2 with the exclusive-or of s of the new input bits, and therefore the size of the evaluated circuit C is $|C| = |C^0| + O(ns)$ gates. The communication overhead is therefore $O(s(|C^0| + ns) + s^2n) = O(s|C^0| + s^2n)$ times the length of the secret-keys (and ciphertexts) used to construct the garbled circuit. (Note that the improved construction in Section 5.2 reduces the size of the new circuit to $|C| = |C^0| + O(\max(n, s))$ and therefore only improves the communication overhead by a constant; the significance of the improvement is with respect to computation.)

The computation overhead is dominated by the oblivious transfers. In Protocol 2 each input bit of P_2 is replaced by s new inputs and therefore $O(ns)$ OTs are required. In Section 5.2 we show how to use only $O(\max(n, s))$ new input bits, and consequently the number of OTs is reduced to $O(\max(n, s))$ (namely $O(1)$ OTs per input bit, assuming $n = \Omega(s)$).

5.1 Efficient Implementation of the Different Primitives

In this section, we describe efficient implementations of the different building blocks of the protocol.

Encryption scheme. Following [20], the construction uses a symmetric key encryption scheme that has indistinguishable encryptions for multiple messages and an elusive efficiently verifiable range. Informally, this means **(1)** that for any two (known) messages x and y , no polynomial-time adversary can distinguish between the encryptions of x and y , and **(2)** that there is a negligible probability that an encryption under one key falls into the range of encryptions under another key, and given a key k it is easy to verify whether a certain ciphertext is in the range of encryptions with k . See [20] for a detailed discussion of these properties, and for examples of easy implementations satisfying them. For example, the encryption scheme could be $E_k(x) = \langle r, f_k(r) \oplus x0^n \rangle$, where k is a pseudo-random function keyed by k , and r is a randomly chosen value.

Commitment schemes. The protocol uses both unconditionally hiding and unconditionally binding commitments. Our goal should be, of course, to use the most efficient implementations of these primitives, and we therefore concentrate on schemes with $O(1)$ communication rounds (all commitment schemes we describe here have only two rounds). Efficient unconditionally hiding commitment schemes can be based on number theoretic assumptions, and use $O(1)$ exponentiations (see, e.g., [13,26]). The most efficient implementation is probably the

one due to Halevi and Micali, which uses a collision-free hashing function and no other cryptographic primitive [14]. Efficient unconditionally binding commitments can be constructed using the scheme of Naor [24], which has two rounds and is based on using a pseudo-random generator.

Oblivious transfer. The protocol needs to use an OT protocol which is secure according to the real/ideal model simulation definition. Candidate protocols can be the protocol of [7] compiled according to the GMW paradigm, or the two-round protocols of [25,2,16] with additional proofs of knowledge.

5.2 Reducing the Number of Oblivious Transfers

Protocol 2 uses a construction which replaces each input bit of P_2 with multiple input bits, providing P_2 with multiple options for encoding each of its inputs. This limits the information that P_1 can gain from corrupting OT inputs (and in particular, P_2 aborts with almost the same probability irrespective of its actual input). The construction replaces each original input wire of P_2 with s new wires, thus increasing the number of input wires of P_2 from n to ns . We show here a probabilistic construction which reduces the number of input wires of P_2 to $\max(4n, 8s)$ (we also describe how to use codes to construct an explicit construction with similar performance). The construction has a direct effect on the overhead of the protocol, since the number of OTs is equal to the number of input wires of P_2 .

We denote the original input bits as w_1, \dots, w_n and the new input bits as w'_1, \dots, w'_m . Our goal is to minimize m . Each w_i is defined as the exclusive-or of a subset of the new input bits. We define the indicator vector z_i as an m -bit binary string whose j th bit is 1 iff w'_j is in the subset of new input bits whose exclusive-or is w_i . The construction described in Protocol 2 corresponds to indicator vectors $z_i = (\underbrace{0 \dots 0}_{(i-1)s} \underbrace{1 \dots 1}_s \underbrace{0 \dots 0}_{(n-i)s})$.

Before analyzing the constructions, let us recall how P_2 encodes its inputs: it chooses random values for the bits w'_1, \dots, w'_m , subject to the constraint that the exclusive-or of any set of new bits which corresponds to an original bit w_i is equal the original value of w_i . P_2 then runs an OT for each of its new input bits. If one of the answers it receives in these OTs is corrupt, it aborts the protocol. Our goal is to make sure that the decision to abort does not reveal information about P_2 's original input (this is the only place that it is used in the proof). It is clear that if P_1 corrupts the inputs of a single OT, then, since each input bit of P_2 is the exclusive-or of several new bits, the decision to abort does not reveal information about any specific input bit of P_2 . This observation must be generalized for the case of P_1 corrupting more OT inputs, and hold with respect to any subset of P_2 's inputs.

Warmup – reusing bits. In order to use less “new” input bits, P_2 must reuse these bits. Assume that P_2 has two input wires w_1, w_2 and that we replace them with $s+1$ new wires, w'_1, \dots, w'_{s+1} . The input values are defined as $w_1 = w'_1 \oplus \dots \oplus w'_s$,

and $w_2 = w'_2 \oplus \dots \oplus w'_{s+1}$ (namely $z_1 = 11 \dots 10$ and $z_2 = 01 \dots 11$). In this case, it is easy to see that any strategy used by a malicious P_1 to corrupt OT values gives it an advantage of at most 2^{-s+1} in identifying a single bit of P_2 's original input (e.g., if P_1 corrupts the '1' inputs of w'_1, \dots, w'_s , then if $w_1 = 1$ P_2 always aborts, whereas if $w_1 = 0$ there is a probability of 2^{-s+1} that P_2 does not abort). However, $w_1 \oplus w_2 = w'_1 \oplus w'_{s+1}$ (namely, $z_1 \oplus z_2 = 10 \dots 01$) and therefore if P_1 corrupts the OT values of both w'_1 and w'_{s+1} it can obtain a non-negligible advantage in learning $w_1 \oplus w_2$. (For example, P_1 can corrupt the '1' inputs of w'_1 and w'_{s+1} . If P_2 does not abort P_1 can conclude that $w'_1 = w'_{s+1} = 0$ and therefore $w_1 = w_2$.)

The attack presented above can be prevented if the exclusive-or of any subset of P_2 's original bits contains at least s new input bits. Namely, if, in the general case, for every non-empty subset $L \subseteq \{1, \dots, n\}$ it holds that the Hamming weight of $\bigoplus_{i \in L} z_i$ is at least s . The two lemmata stated below (which are proved in the full version of the paper) show that this requirement is sufficient to prove that, up to a negligible probability, P_2 's decision to abort is independent of its input values.

Lemma 1. *Suppose that for any set $L = \{i_1, \dots, i_{|L|}\}$ (corresponding to a set $\{w_{i_1}, \dots, w_{i_{|L|}}\}$ of original input wires), the Hamming weight of $z_{i_1} \oplus \dots \oplus z_{i_{|L|}}$ is at least s . Fix the values of any subset of less than s new input wires arbitrarily, and choose the values of all other new input wires uniformly at random. Then for any set $L = \{i_1, \dots, i_{|L|}\}$, the value of the vector $(w_{i_1}, \dots, w_{i_{|L|}})$ is uniformly distributed.*

Lemma 2. *Suppose that for all sets $L = \{i_1, \dots, i_{|L|}\}$ the Hamming weight of $z_{i_1} \oplus \dots \oplus z_{i_{|L|}}$ is at least s . Then, for any two different inputs y and y' of P_2 , the difference between the probability that P_2 aborts the protocol as a result of corrupt OT values when its input is y and when its input is y' is at most 2^{-s+1} .*

Given Lemma 2 it is possible to construct assignments of the new input values to the original input values which ensure that OT corruptions by P_1 do not reveal information about P_2 's input. The constructions are based on ensuring that for any set $S = \{i_1, \dots, i_{|L|}\}$ the Hamming weight of $z_{i_1} \oplus \dots \oplus z_{i_{|L|}}$ is at least s . We describe below a randomized construction which achieves this property. As was pointed to us by David Woodruff, an explicit construction can be achieved using any explicit linear code from $\{0, 1\}^s$ to $\{0, 1\}^{O(s)}$, for which any two codewords have a distance of at least $\Omega(s)$ (Justesen codes are an example of such a code).

The randomized construction. We define $4n$ new input bits for P_2 . Assume, without loss of generality, that $n > 2s$. (Otherwise add dummy input bits. Therefore the exact number of new input bits is $\max(4n, 8s)$.) The mapping between the n old input bits and the $4n$ new input bits is chosen randomly in the following way: each original input bit w_i is defined to be equal to the exclusive-or of a uniformly chosen subset of the new input bits (in other words, z_i is a uniformly distributed string of $4n$ bits).

We examine the probability that there is a subset $L \subseteq \{0, 1\}^n$ for which the Hamming weight of $\bigoplus_{i \in L} z_i$ is less than s : Consider any subset L , then $\bigoplus_{i \in L} z_i$ is a uniformly distributed string with $4n > 8s$ bits, with an expected Hamming weight of $2n$. Let X_j be a random variable which is set to 1 if the j th bit in this string is 1. Note that $s/4n < 1/8$ by our assumption that $n > 2s$. We have:

$$\Pr \left[\sum_{j=1}^{4n} X_j < s \right] = \Pr \left[\frac{\sum X_j}{4n} < \frac{s}{4n} \right] < \Pr \left[\frac{\sum X_j}{4n} < \frac{1}{8} \right] \leq \Pr \left[\left| \frac{\sum X_j}{4n} - \frac{1}{2} \right| > \frac{3}{8} \right]$$

Applying the Chernoff bound, we have that $\Pr \left[\sum_{j=1}^{4n} X_j < s \right] = < 2e^{-\frac{(3/8)^2}{2(1/2)(1/2)}4n} = 2e^{-9n/8}$. There are a total of 2^n subsets of the original input bits, and therefore the probability that any of them is equal to the exclusive-or of less than s new input bits is bounded by $2^n 2e^{-9n/8} \approx 2^{(1-9/8 \log(e))n} \approx 2^{-0.6n} < 2^{-1.2s}$. Lemma 2 therefore implies that with probability $1 - 2^{-1.2s}$ the construction suffices for our proof of security.

Choosing the strings z_i . In order to use the above construction, the parties must construct a circuit that has $4n$ new input bits for P_2 . Furthermore, the parties must define n random strings z_i of length $4n$ and then have the circuit map P_2 's i^{th} input bit according to the string z_i (as described above). This can be done in two ways. One possibility is to choose the mapping once and for all and hardwire it into the protocol specification. This is problematic because then there is a negligible probability that the protocol is not secure (in any execution). Thus, the mapping should instead be chosen as part of the protocol execution (because negligible failure in any execution is allowed). Fortunately, P_2 can singlehandedly choose the strings z_1, \dots, z_n in the first step of the protocol and send them to P_1 . The reason why this is fine is because this entire issue only arises in the proof of the case that P_1 is corrupted (indeed, for the case of a corrupted P_2 there is no need to split P_2 's input bits at all).

Acknowledgments

We would like to thank Yuval Ishai, Moni Naor, Adam Smith and David Woodruff for helpful discussions about this work.

References

1. G. Aggarwal, N. Mishra, and B. Pinkas. Secure Computation of the k -th Ranked Element. In *EUROCRYPT 2004*, Springer-Verlag (LNCS 3027), 40–55, 2004.
2. B. Aiello, Y. Ishai, and O. Reingold. Priced Oblivious Transfer: How to Sell Digital Goods. In *EUROCRYPT 2001*, Springer-Verlag (LNCS 2045), 119–135, 2001.
3. B. Barak and Y. Lindell. Strict Polynomial-Time in Simulation and Extraction. *SIAM Journal on Computing*, 33(4):783–818, 2004.
4. D. Beaver. Foundations of Secure Interactive Computing. In *CRYPTO'91*, Springer-Verlag (LNCS 576), pages 377–391, 1991.

5. R. Canetti. Security and Composition of Multiparty Cryptographic Protocols. *Journal of Cryptology*, 13(1):143–202, 2000.
6. R. Cramer, I. Damgard, and B. Schoenmakers. Proofs of Partial Knowledge and Simplified Design of Witness Hiding Protocols. In *CRYPTO'94*, Springer-Verlag (LNCS 839), pages 174–187, 1994.
7. S. Even, O. Goldreich and A. Lempel. A Randomized Protocol for Signing Contracts. In *Communications of the ACM*, 28(6):637–647, 1985.
8. O. Goldreich. *Foundations of Cryptography: Volume 1 – Basic Tools*. Cambridge Univ. Press, 2001.
9. O. Goldreich. *Foundations of Cryptography: Volume 2 – Basic Applications*. Cambridge Univ. Press, 2004.
10. O. Goldreich and A. Kahan. How To Construct Constant-Round Zero-Knowledge Proof Systems for NP. *Journal of Cryptology*, 9(3):167–190, 1996.
11. O. Goldreich, S. Micali and A. Wigderson. How to Play any Mental Game – A Completeness Theorem for Protocols with Honest Majority. In *19th STOC*, pages 218–229, 1987. For details see [9].
12. S. Goldwasser and L. Levin. Fair Computation of General Functions in Presence of Immoral Majority. In *CRYPTO'90*, Springer-Verlag (LNCS 537), pages 77–93, 1990.
13. S. Goldwasser, S. Micali and R.L. Rivest, A Digital Signature Scheme Secure Against Adaptive Chosen-Message Attacks. *SIAM J. Comput.* 17(2): 281–308, 1988.
14. S. Halevi and S. Micali, Practical and Provably-Secure Commitment Schemes from Collision-Free Hashing, *CRYPTO 1996*, Springer-Verlag (LNCS 1109), pages 201–215, 1996.
15. S. Jarecki and V. Shmatikov. Efficient Two-Party Secure Computation on Committed Inputs. In these proceedings (Eurocrypt '2007).
16. Y.T. Kalai. Smooth Projective Hashing and Two-Message Oblivious Transfer. In *EUROCRYPT 2005*, Springer-Verlag (LNCS 3494), pages 78–95, 2005.
17. J. Katz and Y. Lindell. Handling Expected Polynomial-Time Strategies in Simulation-Based Security Proofs. In the *2nd Theory of Cryptography Conference (TCC)*, Springer-Verlag (LNCS 3378), pp. 128–149, 2005.
18. J. Kilian. Founding Cryptography on Oblivious Transfer. In *20th STOC*, pages 20–31, 1988.
19. M. Kiraz and B. Schoenmakers. A Protocol Issue for the Malicious Case of Yao's Garbled Circuit Construction. In *Proceedings of 27th Symposium on Information Theory in the Benelux*, 283–290, 2006.
20. Y. Lindell and B. Pinkas. A Proof of Yao's Protocol for Secure Two-Party Computation. To appear in the *Journal of Cryptology*. Also appeared as *Cryptology ePrint Archive*, Report 2004/175, 2004.
21. D. Malkhi, N. Nisan, B. Pinkas and Y. Sella. Fairplay – A Secure Two-Party Computation System. In the *13th USENIX Security Symposium*, pages 287–302, 2004.
22. S. Micali and P. Rogaway. Secure Computation. Unpublished manuscript, 1992. Preliminary version in *CRYPTO'91*, Springer-Verlag (LNCS 576), pages 392–404, 1991.
23. P. Mohassel and M.K. Franklin. Efficiency Tradeoffs for Malicious Two-Party Computation. In the *9th PKC conference*, Springer-Verlag (LNCS 3958), pages 458–473, 2006.
24. M. Naor. Bit Commitment Using Pseudorandomness. *Journal of Cryptology*, 4(2):151–158, 1991.

25. M. Naor and B. Pinkas. Efficient Oblivious Transfer Protocols. In the *12th SODA*, pages 448-457, 2001.
26. T.P. Pedersen. Non-Interactive and Information-Theoretic Secure Verifiable Secret Sharing. In *CRYPTO'91*, Springer-Verlag (LNCS 576), pages 129-140, 1992.
27. M. Rabin. How to Exchange Secrets by Oblivious Transfer. Tech. Memo TR-81, Aiken Computation Laboratory, Harvard U., 1981.
28. D. Woodruff. Revisiting the Efficiency of Malicious Two-Party Computation. In these proceedings (Eurocrypt '2007).
29. A. Yao. How to Generate and Exchange Secrets. In *27th FOCS*, pages 162-167, 1986.