

Fast and Maliciously Secure Two-Party Computation Using the GPU (Full version)*

Tore Kasper Frederiksen¹ and Jesper Buus Nielsen¹

Department of Computer Science, Aarhus University
jot2re@cs.au.dk, jbn@cs.au.dk [†]

Abstract. We describe, and implement, a maliciously secure protocol for two-party computation in a parallel computational model. The protocol is based on cut-and-choose of Yao’s garbled circuit and an efficient oblivious transfer extension. The implementation is done using CUDA and yields fast results in a financially feasible and practical setting by using a consumer grade CPU and GPU. Our protocol introduces a novel construction in order to verify consistency of the garbled circuit constructor’s input in a parallel and maliciously secure setting.

1 Introduction

Secure two-party computation (2PC) is the area of cryptography concerned with two mutually distrusting parties who wish to securely compute an arbitrary function on their private input in order for each of them to learn some private output. The area was introduced in 1982 by Andrew Yao [Yao82], specifically for the *semi honest* case where both parties are assumed to follow the prescribed protocol. Yao showed how to construct such a protocol using a technique referred to as the *garbled circuit approach*. Later, a solution in the *malicious* setting, where one of the parties might deviate from the prescribed protocol in an arbitrary manner, was given in [GMW87]. Unfortunately this protocol was very inefficient as it depended on public-key operations for each Boolean gate in the circuit describing the function to compute. However, much research has since been done in the area of 2PC, resulting in practically efficient protocols secure both against both semi honest and malicious adversaries [LP07, LP11, NO09, NNOB12, HEKM11, PSSW09].

In general, the protocols secure against a semi honest adversary can become secure against a malicious adversary by adding an additional layer of security. This can be either by compiling a semi honest secure protocol to a maliciously secure protocol [GMW87] or using the cut-and-choose approach where several instances of a semi honest secure protocol are executed in parallel with some random instances being completely revealed to verify that the other party has behaved honestly. However, novel approaches to achieve malicious security do exist, such as the idea of MPC-in-the-head from [IPS08, LOP11] or by embedding the cut-and-choose part at a lower level of the protocol as done in [NO09] or [NNOB12]. However, assuming access to a large computer grid and using the cut-and-choose approach in a parallel manner has yielded slightly better results than [NNOB12] as described in [KSS12].

Motivation. The area of 2PC and multi-party computation (MPC) (when more than two parties supply input) is very interesting as efficient solutions yield several practical applications. The first case of this is described in [BCD⁺09] where MPC was used for deciding the market clearing price of sugar beets in Denmark. Several other applications for 2PC and MPC includes voting, anonymous identification, privacy preserving database queries etc. For this reason we believe that it is highly relevant to find practically

* This is a revision handling a bug in the way we previously avoided selective failure attacks on the input to the oblivious transfers in our protocol. Notice that this bug was discovered *after* the extended abstract version of this work [FN13] was published.

[†] Partially supported by the Danish Council for Independent Research via DFF Starting Grant 10-081612. Partially supported by the European Research Commission Starting Grant 279447.

efficient protocols for 2PC and MPC. Most previous approaches have focused on doing this in a sequential model [NNOB12, LP07, LP11]. However, considering the recent evolution of processors we see that the speed of a processor seems to converge, whereas the amount of cores in a processor seems to increase. This in turn implies that the increase in processing power in the future will come from having many cores working in parallel. Thus, constructing both algorithms and cryptographic protocols that work well in a parallel model will be paramount for hardware based efficiency increases in the future. For this reason we have chosen to take a parallel approach to increase the practical speed of 2PC. Previous work taking the parallel approach for efficient implementations of MPC starts with [PDL11] where a cluster of either CPUs or GPUs is used to execute 3072 semi honest protocols for 1-out-of-2 oblivious transfer (OT) followed by gate garbling/degarbbling in parallel¹. In [KSS12] the authors use up to 512 cores of the Ranger cluster in the Texas Advanced Computing Center to do OTs along with circuit garbling/degarbbling in parallel to achieve malicious security using the cut-and-choose approach. In this manner they manage to use the inherent parallelism of the cut-and-choose approach to achieve very fast and maliciously secure 2PCs. Any other work taking a parallel approach to cryptography that we know of focuses either on attacks [XLZ11] or simultaneous applications of more primitive cryptographic computations [NIK12].

Contributions. Our main contributions are a new approach to ensure consistency of the garbler’s inputs in cut-and-choose based protocols of Yao’s garbled circuit, along with a careful implementation of a protocol using this approach implemented on a *Same Instruction, Multiple Data (SIMD)*, or *Parallel Random Access Model (PRAM)* computation device. More specifically the protocol we implement is secure in the *Random Oracle Model (ROM)*, *OT-hybrid* model and loosely based on the protocol of [LP07] but combined with several newer optimizations along with the *OT extension* of [NNOB12].² Computationally our protocol relies solely on symmetric primitives, except for a few seed OTs used in the OT extension which only need to be done once for each pair of parties. Furthermore, our protocol is of constant round complexity and, assuming access to enough cores, computationally bounded only by the number of layers in the circuit to be computed and the block size of a hash function. Using a NVIDIA GPU as our SIMD device, we make several experiments and we show that this approach is one of the fastest documented assuming a “practical”, yet malicious, setting.³

Notation. We let \parallel denote string concatenation and let $r[i]$ be the i ’th element of a string r . We let ℓ be a statistical security parameter such that the probability of an unbounded adversary compromising the security of our protocol is at most $2^{-\ell}$. We then let κ be the computational security parameter such that the probability of a probabilistic polynomial time bounded adversary compromising the security of our protocol is negligible in κ . We then let $H(\cdot)$ denote a hash function with a digest of κ bits (in our implementation this will be 160 bits) and block size ρ (in our implementation this will be 448 bits). We assume that Alice is the circuit *constructor* and Bob is the circuit *evaluator* and we will use their names and roles interchangeably.

Overview. The rest of the paper is organized as follows: We start in Section 2 with an introduction to the idea of parallel implementations and the overall structure of our computation device of choice; the GPU. Then in Section 3 we go through the overall structure of our protocol. In Section 4 we go through the ideas we used to make it suitable for the SIMD model. In Section 5 we discuss the security and complexity of our approach. This is followed by Section 6 where we discuss the implementation details. Then in Section 7 we review our results and end up with a discussion of future work in Section 8.

¹ 1-out-of-2 OT is the protocol where the first party, Alice, gives as input two bitstrings, and the second party, Bob, gives a single bit. If Bob’s bit is 0 then he learns Alice’s first bitstring, if it is 1 then he learns Alice’s second bitstring. However, Bob never learns more than exactly one of these bitstrings and Alice does not learn anything.

² An OT extension is a protocol for transforming a small number of OTs into a large number of OTs. Such an extension often makes sense as it can be done using only cheap cryptographic primitives such as hash functions. Whereas “plain”, non-extended OTs, generally rely on public-key cryptography which is significantly slower than computing digests of a hash function.

³ We refer to “practical” as either financially feasible using consumer hardware and/or having a liberal statistical security parameter.

2 Background

Parallel Approach. In our approach we assume access to a massive parallel computation device which is capable of executing the same instruction on each processor in parallel, but on different pieces of data. This is in a sense the simplest way of modeling parallel computation, as a device capable of executing distinct instructions on distinct pieces of data is clearly also capable of executing the same instruction on distinct pieces of data. Furthermore, our protocol does not make any assumption on whether such a device has access to shared memory between the processors, or only access to local memory. This applies completely for write privileges, but also for read privileges with only a constant memory usage penalty if shared memory is not available.

GPGPU. We decided to implement our protocol using the GPU, the motivation being that GPUs are part of practically all mid- to high-end consumer computers. Furthermore, using the GPU eliminates the security problems of outsourcing the computation to a non-local cluster. Also, assuming access to a local cluster seems to be an unrealistic assumption for many practical applications. Finally, using gaming consoles or multi-cores CPUs might also have been an option. However, even the latest and best of these have orders of magnitude cores less than the latest GPUs.

CUDA. Our implementation is done using the CUDA framework which is an extension to C and C++ that allows using NVIDIA GPUs for general computational tasks. This is done by making CUDA programs. Such a program does not purely run on the GPU. It consists of both general C classes, which run on the CPU, and CUDA classes which run on the GPU.

In order to motivate our specific implementation choices it is necessary to describe a general CUDA enabled GPU [KH10]: Each GPU consists of several (up to 192) *streaming multiprocessors* (SM), each of these again contains between 8 and 192 *streaming processors* (SP), depending on the architecture of the GPU. Each of the SPs within a given SM always performs the same operations at a given point in time, but on different pieces of data. Furthermore, each of these SMs contains 64 KB of *shared memory* along with a few kilobytes of constant cache, which all of the SPs within the given SM must share. For storage of variables each SM contains 65.536 32-bit registers which are shared amongst all the SPs. Thus all the threads being executed by a given SM must share all these resources.

We now introduce some notation and concepts which are used in the general purpose GPU community and which we will also use in this paper; a GPU is called a *device* and the non-GPU parts of a computer is called the *host*. This means that the CPU, RAM, hard drive etc., are part of the host. The code written for the host will be used to interact with the OS, that is, it will do all the IO operations needed by the CUDA program. The host code is also responsible for copying the data to and from the device, along with launching code on the device. Each procedure running on a device is called a *kernel*. Before launching a kernel the host code should complete all needed IO and copy all the data needed by the kernel to the device's RAM. The RAM of the device is referred to as *global memory*. After a kernel has terminated the host can copy the results from the global memory of the device to its own memory, before it launches another kernel.

A kernel is more than just a procedure of code, it also contains specifications of how many times in parallel the code should be executed and any type of synchronization needed between the parallel executions. A kernel consists of code which is executed in a *grid*. A grid is a 2-dimensional matrix of *blocks*. Each block is a 3-dimensional matrix of threads. Each thread is executed once and takes up one SP during its execution. When all the threads, of all the blocks in the grid, have been executed the kernel terminates. The threads in each block are executed in *warps*, which is a sequence of 32 threads. Thus, the threads must be partitioned into blocks in multiples of the warp size, and contain no branching. The threads can then be executed completely independent and in arbitrary order.

Furthermore, to achieve the fastest execution time one should *coalesce* the data in the global memory. That is, to “sort” the data such that the word thread 1 needs is located next to the word thread 2 needs and so on. This makes it possible to load these 32 words for the warp in one go, thus limiting the usage of bandwidth, and in turn significantly increasing the speed of the program. This advice on memory organisation is also relevant for the data in the shared memory. Finally, it is a well known fact [Cor12] that the bottleneck

for most applications of the massive parallelism offered by CUDA is the memory bandwidth, thus it should always be a goal to limit the frequency of which a program access data in the global memory.

Maliciously Secure Garbled Circuits.

Generic Garbled Circuits. For completeness we now sketch how a generic garbled circuit is constructed. We are given a Boolean circuit description, C , of the Boolean function we wish to compute, f , from which we construct a garbled circuit, GC . For simplicity we assume that each gate consists of two input wires and one output wire. However, we allow the output wire to split into two or more if the output of a given gate is needed as input to more than one other gate. Each wire in C has a unique label, and we give the corresponding wire in GC the same label. Each wire w has two keys associated, k_w^0 and k_w^1 , which are independent uniformly random bitstrings. Here k_w^0 represents the bit 0 and k_w^1 represents the bit 1. If the bit on wire w in C is 0, then the value on wire w in GC will be k_w^0 , otherwise it will be k_w^1 . Each gate in GC consists of a *garbled computation table*. This table is used to find the correct value of the output wire given the correct keys for the input wires. For a gate of C call the left input wire l , the right input wire r and the output wire o . Assume the functionality of the gate is given by $G(\sigma, v) = \rho$ where $\sigma, v, \rho \in \{0, 1\}$, then the garbled computation table is a random permutation of the four ciphertexts $C_{\sigma, v} = E_{k_l^\sigma}(E_{k_r^v}(k_o^\rho)) = E_{k_l^\sigma}(E_{k_r^v}(k_o^{G(\sigma, v)}))$ for all four possible input pairs, (σ, v) , using some symmetric encryption function, $E_{\text{key}}(\cdot)$. I.e., the entries in the garbled computation table consists of “double encryptions” of the output wire’s values, where the keys for each double encryption corresponds to exactly one combination of the input wires’ values. The encryption algorithm is constructed such that given k_l^σ and k_r^v it is possible to recognize and correctly decrypt $C_{\sigma, v}$, but it is not possible to learn any information about the remaining three encryptions.

Now let Alice’s input be denoted by x and Bob’s input be denoted by y . Then the generic version of semi honestly secure 2PC based on GCs [LP09] goes as follows:

1. Alice starts by constructing a GC, GC , for a Boolean circuit, C , computing the desired function, $f(x, y) = (f_1(x, y), f_2(x, y))$ where Alice is supposed to learn $f_1(x, y)$ and Bob is supposed to learn $f_2(x, y)$.
2. Alice sends to Bob the garbled computation table. She also sends both the keys of the output wires of Bob’s output, along with the bit each of these keys represent.
3. Alice now sends keys for the first $|x|$ input wires, corresponding to her desired input. That is, for $i = 1, \dots, |x|$ she sends either k_i^0 or k_i^1 corresponding to her i ’th input bit.
4. Next Alice and Bob complete a 1-out-of-2 OT protocol $|y|$ times:
 - (a) For $i = |x| + 1, \dots, |x| + |y|$ Alice obliviously sends the keys k_i^0 and k_i^1 to Bob.
 - (b) Bob then chooses exactly one of these keys for each i , without Alice knowing which one.
5. Now Bob has the circuit along with a set of input keys. However, he does not know whether each of the keys to Alice’s input represents a 0 or 1 bit, but he does know that the keys for his input represents the correct bits in accordance with his bitstring, y .
6. Bob now degarbles the circuit and learns output keys for all the output wires. He then compares these with the output keys he got from Alice and finds his output bits. He then sends to Alice the keys for the output wires corresponding to her output.
7. Using these keys Alice finds the bits they represent and thus her output.

As soon as any of the players deviate from the protocol, this scheme breaks down completely.

Optimized Garbled Circuits. Several universal techniques for optimizing the generic (both malicious and semi honest security) garbled circuit approach for 2PC exist, which we also implement in our protocol. We now go through the specific and optimized construction of garbled circuits we use.

First notice that we did not specify exactly how to determine which entry in the garbled computation table is the correct one to decrypt given two input wire keys. This is because several different approaches for this exist. However, one efficient approach is the usage of permutation bits [NPS99]. The idea is to associate a single *permutation bit*, $\pi_i \in \{0, 1\}$, with each wire, i , in the circuit. The value on this wire is then defined as $k_i^b \parallel c_i$ where $c_i = \pi_i \oplus b$ with b being the bit the wire should represent. We call c_i the *external value*.

The entries in the garbled computation table are then sorted according to the external values. That is, if the external value on both the left and right wire is 0 then the key in the *first* entry of the table will be encrypted under these two wire keys. If instead the external value on the right wire is 1, then the key in the *second* entry of the table will be encrypted under the left and right wire keys. Thus, we view the external values as a binary number, specifying an entry in the garbled computation table. More formally, entry c_l, c_r of the garbled computation table is specified as follows:

$$c_l, c_r : E_{k_l^{b_l}, k_r^{b_r}}^{Gid \| c_l \| c_r} \left(k_o^{G(b_l, b_r)} \| c_o \right),$$

where $c_l = \pi_l \oplus b_l$, $c_r = \pi_r \oplus b_r$, $c_o = \pi_o \oplus G(b_l, b_r)$ and Gid is a unique gate ID. This means that given the keys of the input wires the evaluator can decide exactly which entry he needs to decrypt, without learning anything about the bits the input wires represent. Next, see that the encryption function for the keys in the garbled computation tables can be described as follows:

$$E_{k_l, k_r}^s(k_o) = k_o \oplus \text{KDF}^{|k_o|}(k_l, k_r, s),$$

where $\text{KDF}^{|k_o|}(k_l, k_r, s)$ is a *key derivation function* with an output of $|k_o|$ bits, independent of the two input keys, k_l and k_r in isolation, and which depends on the value of some salt, s . If we assume the ROM, then we are able to specify the KDF as follows [PSSW09]:

$$\text{KDF}^{|k_o|}(k_l, k_r, s) = H(k_l \| k_r \| s).$$

This means that the encryption function can be reduced to a single invocation of a robust hash function with output length κ (assuming $\kappa \geq |k_o|$) along with an XOR operation.

Next we describe the optimization from [KS08] which will make it possible to evaluate all the XOR gates in the circuit for “free”. Free here means that no garbled computation table needs to be constructed or transmitted, and no encryption needs to be done in order to evaluate such a gate. The only thing we need to do to get this possibility is to put a constraint on the way the wire keys are constructed. The constraint is very simple, assuming that we wish to construct the keys for wire i , then it must be the case that

$$k_i^1 = k_i^0 \oplus \Delta,$$

where Δ is a *global key*, used in the keys for all wires in the GC. Regarding the external values, this implies the following:

$$\pi_i \oplus 1 = \pi_i \oplus 0 \oplus 1.$$

So, in order to compute an XOR gate simply compute XOR of the keys of the two input wires of the gate, that is:⁴

$$k_o \| c_o = k_l \oplus k_r \| c_l \oplus c_r.$$

Finally, see that a row of the garbled computation table can be eliminated using the approach of [NPS99]. To see this remember that a single Boolean gate takes up four entries of κ bits. However, when using a KDF for encryption we can simply define one of the output keys to be the result of this KDF on one input key pair. This key pair is the one where the external values are 0, i.e. $c_l = 0$ and $c_r = 0$. In short, we must define one of the output keys, and an external value as follows:

$$k_o^{G(\pi_l \oplus 0, \pi_r \oplus 0)} \| c_o = \text{KDF}^{\kappa+1} \left(k_l^{\pi_l \oplus 0}, k_r^{\pi_r \oplus 0}, \text{Gid} \| 0 \| 0 \right).$$

Depending on the type of gate, this again uniquely specifies the permutation bit of the output wire from this calculation:

$$c_o = \pi_o \oplus G(\pi_l \oplus 0, \pi_r \oplus 0).$$

The three remaining entries in the garbled computation table are then the appropriate encryptions of the output key or the output key XOR the global difference Δ .

⁴ Notice that when using the free-XOR approach Alice cannot simply send Bob both keys on the output wires since this would make him learn Δ which can be used to extract information about Alice’s input. Several solutions to this exist, for example, the final garbled gates can simply encrypt the plain bits of their output instead of new keys.

Optimized Approaches to Cut-and-Choose Malicious Security. In general OT is an expensive primitive, and if the evaluator has a large input to the circuit this can contribute significantly to the execution time of the whole protocol. However, the amount of “actual” OTs we need to complete can be significantly reduced by using an OT extension: Beaver showed in [Bea96] that given a number of OTs it is possible to “extend” these to give a polynomial number of random OTs which can easily be changed to specific OTs. Thus, making it possible to do a few OTs once, and extend these almost indefinitely. The idea of an OT extension has been optimized even further in [IKNP03] and [NNOB12] to yield significant practical advantages. Our protocol uses a slightly modified version of the OT extension presented in [NNOB12].

Even when using a maliciously secure OT extension the cut-and-choose approach in itself is unfortunately not enough to make a semi honestly secure protocol maliciously secure efficiently. In fact, several problems arise from using cut-and-choose to get security against a malicious adversary, these problems can be categorized as follows:

1. “Consistency of input bits”; both parties need to use the same input in all the cut-and-choose instances to ensure that the majority of the garbled circuit evaluations are consistent and that a corrupt evaluator does not learn the output of the function on different inputs.
2. “Selective failure attack”; we must make sure that both the keys the constructor inputs in the OT phase are correct, to avoid giving away a particular bit value of the evaluator’s input, depending on failure or not of the evaluation.

Consistency of the evaluator’s input is achieved by simply doing one OT for each of his input bits and have the constructor input the 0-keys, respectively 1-keys, for each of his input wires in *all* the garbled circuits. In this manner the evaluator only gets to choose his input bits once during the OTs and thus it will be impossible for him to be inconsistent between the different garbled circuits.

Consistency of the Constructor’s Input. Letting $|x|$ be the size of the constructor’s input and ℓ the statistical security parameter then the first problem can be solved using $O(|x| \cdot \ell^2)$ commitments to verify consistency in all possible cut-and-choose cases [LP07]. A more efficient approach is to construct a Diffie-Hellman pseudo random synthesizer, which limits the complexity to $O(|x| \cdot \ell)$ symmetric and public-key operations and also solves the selective failure attack [LP11]. Yet another solution is based on claw-free functions [SS11]. Our solution is different; we solve the problem by constructing a circuit extension making it possible to verify, with statistical security, that the constructor has been consistent in her inputs to the circuit. Assume the function we wish to compute is defined by f as $f(x, y) = z$. We then define a new function f' as $f'(x', y') = f'((x||s), y||r) = z||t$ where $s \in_R \{0, 1\}^\ell$, $r \in_R \{0, 1\}^{|x|+\ell-1}$ and $t \in \{0, 1\}^\ell$. To compute t we define a matrix $M_1 \in \{0, 1\}^\ell \times \{0, 1\}^{|x|}$ where the i 'th row is the first $|x|$ bits of $r \ll i$ where \ll denotes the bitwise left shift. Specifically the j 'th bit of the i 'th row is the $i + j$ 'th bit of the binary vector r , i.e. $M_1[i, j] = r[i + j]$. Using this matrix the computation of t is defined as $t = (M_1 \cdot x) \oplus s$, assuming all binary vectors are in column form.

With this modification the new function computes the same as the original, but requires ℓ extra random bits of input from the constructor and $|x| + \ell - 1$ extra random bits from the evaluator. However, the new function returns ℓ extra bits to the evaluator. These ℓ extra bits will work as digest bits and can be used to check that the constructor is consistent with her inputs to the GCs by verifying that they are the same in all the garbled circuits which are evaluated. However, it must clearly still be the case that honest parties input the same input $x||s$ and $y||r$ for each of the garbled circuits.

This augmentation works since the new function computes, besides the original functionality, a family of universal hash functions where the auxiliary input from both parties defines a particular hash function from this family. The auxiliary output of the augmented function is then the digest of the constructor’s input in this universal hash function. The proof that the augmentation is indeed a family of universal hash functions was shown in [MNT90]. Intuitively, if the constructor’s tries to give inconsistent inputs between the different garbled circuits, then this will result in different digests in each of these garbled circuits, except with probability $2^{-\ell}$. Thus, the evaluator will be able to detect if the constructor has given inconsistent input. Furthermore, notice that this augmentation does not give Bob any information about Alice’s true

input because Alice gets to mask each bit of $M_1 \cdot x$ by a uniformly random bit selected by herself, i.e. the s vector.

Selective Failure. In [SS11] the problem of a selective failure attack is solved using a special version OT, known as *committing OT*. In [LP07] it is shown how to do this using a circuit extension which increases the amount of input bits of the evaluator from $|y'|$ to $\max(4 \cdot |y'|, 8 \cdot \ell)$. This is the solution we are using in our protocol. More specifically, what we do is to choose a random binary matrix $M_2 \in_R \{0, 1\}^{\max(4 \cdot |y'|, 8 \cdot \ell)} \times \{0, 1\}^{|y'|}$ and a random binary vector $y'' \in_R \{0, 1\}^{\max(4 \cdot |y'|, 8 \cdot \ell)}$ but under the constraint that $M_2 \cdot y'' = y'$ where y' is the “true” binary input of the evaluator to the, in our case already augmented, functionality f' . Thus this new functionality computes exactly the same as the original. Still, the idea of this approach is that if a selective failure attack is done, using the augmented function will not leak any useful information as the entire vector y'' is random learning a single bit of this vector will only give the adversary a negligible advantage in learning one of the constructor’s true input bits. This follows from the fact that the other bits of y'' will be used to hide each of the actual bits of y' . The details and a full proof of security of this approach can be found in [LP07].

3 Protocol Description

We now describe the overall structure of our protocol. For simplicity we assume that only the evaluator is supposed to receive output from the computation. If we wish to compute a circuit where the constructor should also receive output then the circuit extension approach of [LP07], or the signed output approach of [SS11], will work directly in our protocol and be scalable in parallel.

Abstractly our protocol can be described as follows:

1. The parties agree on a statistical security parameter, ℓ , such that the probability of a total breakdown is at most $2^{-\ell}$. They then agree on a binary function, $f(x, y) = z$ where the constructor inputs x , the evaluator inputs y and the evaluator learns the result z .
2. This functionality is extended to enforce consistency on the constructor’s input. The constructor extends her input x to x' such that $x' = x \| s$ where $s \in_R \{0, 1\}^\ell$, similarly the evaluator extends his input y to y' such that $y' = y \| r$ where $r \in_R \{0, 1\}^{|x|+\ell-1}$. We then define a new functionality as $f'(x', y') = f(x, y) \| t$ where the value t is computed as $(M_1 \cdot x) \oplus s$. Finally, M_1 is a $\ell \times |x|$ binary matrix where the i ’th row is the first $|x|$ bits of r shifted i bits to the left.
3. The evaluator then decides on a random matrix $M_2 \in_R \{0, 1\}^{|y''|} \times \{0, 1\}^{|y'|}$ and a random input y'' such that $M_2 \cdot y'' = y'$ where we let $|y''| = \max(4 \cdot |y'|, 8 \cdot \ell)$. He sends this matrix to the constructor and they both agree on a Boolean function, f'' with this M_2 embedded such that $f''(x', y'') = f'(x', M_2 \cdot y'') = f'(x', y') = f(x, y) \| t$. Thus the Boolean functionality to be computed is extended once again.
4. The parties then agree on a binary circuit computing the functionality $f''(x', y'')$ which we call C . The constructor constructs $\ell' = 3.22 \cdot \ell$ GCs in parallel.⁵
5. The constructor and evaluator engage in OT in order for the evaluator to learn the keys corresponding to his input for all ℓ' circuits. We call this the *OT phase*.
 - (a) The constructor and evaluator complete an OT extension which is 1-out-of-2 OTs of random bitstrings.
 - (b) For each of these OTs the constructor extends the two random outputs to two $\ell' \cdot \kappa$ “random” bitstrings. The first representing the 0-keys of the ℓ' garbled circuits and the other the 1-keys. This is done by using the inputs to the OT as seeds for a hash function.
 - (c) Similarly the evaluator extends his output of each OT to a $\ell' \cdot \kappa$ “random” bitstring, representing either the 0 or 1 keys of the ℓ' garbled circuits depending on his choice in the OT.
 - (d) From the circuit generation the constructor will have a 0 and 1 key for given wire in each GC. The constructor then XORs each of the “random” bitstrings she learned from the OT extension with the appropriate keys from the circuit generation and sends all these differences to the evaluator.

⁵ The constant increase in the amount of GCs stems from the fact that cut-and-choose of ℓ circuits only corresponds to statistical security of $2^{-0.311\ell}$ [LP11].

- (e) The evaluator uses these bitstrings to find the correct input keys for the GCs by a simple XOR operation.
- 6. The constructor then commits to each of the ℓ' GCs along with the input keys by hashing them and sending the digests to the evaluator. These digests make it possible to avoid sending half of the garbled computation tables as mentioned in [GMS08].
- 7. The constructor then commits to her input bits by computing a hash digest, concatenated with some random salt, of the keys to each of the garbled circuits in correspondance with her input bits x' to the function f'' . She sends these digests to the evaluator.
- 8. The parties then select $\ell'/2$ circuits for verification using a coin-tossing protocol⁶ and the constructor sends the random seeds used to generate these circuits to the evaluator. We call this and the following three steps for the *cut-and-choose phase*. The coin-tossing can for example be realized as follows:
 - (a) The evaluator chooses a random string $\rho_2 \in_R \{0, 1\}^\ell$ and commits to this by sending a hash digest of it to the constructor.
 - (b) The constructor also chooses a random string, $\rho_1 \in_R \{0, 1\}^\ell$ and sends this to the evaluator.
 - (c) The evaluator sends the string ρ_2 to the constructor, who computes a hash digest and verifies that it matches the digest sent in the first step of the coin-tossing subprotocol.
 - (d) The constructor and evaluator now defines $\rho = \rho_1 \oplus \rho_2$.
 - (e) The parties use the string ρ to deterministically select a subset of ℓ' with size $\ell'/2$.
- 9. Using the seeds the evaluator regenerates the garbled circuits' garbled computation tables along with the input keys and verifies that they are correct by hashing them and checking equality with the digests he received in Step 6.
- 10. After this check the constructor sends the input keys in correspondance with her input, the random salt she used in Step 7 along with the garbled computation tables of the $\ell'/2$ circuits for which the evaluator was *not* given the seeds.
- 11. The evaluator then hashes the garbled computation tables of these circuits and verifies them against the hash digests he received in Step 6. He then hashes the constructor's input keys concatenated with the salt and verifies the digests against the values he received in Step 7. Finally, he evaluates the circuits to achieve their outputs. He then checks that the digest bits are the same in all evaluation circuits. We call this the *evaluation phase*.
- 12. If all checks above pass he takes the majority of the decrypted outputs of the $\ell'/2$ circuits to be the overall output of the protocol, otherwise he aborts.

4 Specific Details

We now give a more technical description of the abstract parts of the protocol.

The Garbled Circuit. We construct the GCs in very much the same manner as described in [Yao82]. However, we use the optimizations for free XOR [KS08], garbled row reduction [NPS99] along with an efficient encryption function using permutation bits [PSSW09]. The overall protocol for achieving malicious security is very similar to that of [LP07], but combined with the sending of circuit seeds [GMS08] and the efficient OT extension of [NNOB12]. Our optimizations arrive in the way we construct and evaluate the GCs in a scalable parallel manner, along with our approach to verify consistency of the constructor's input.

Constructing the Garbled Circuits. We turn the augmented function f'' , into a circuit description which we then parse. The parsing consists of finding all the gates which can be computed using only the input wires, calling this set of gates for layer 0. We then find all the gates, not in layer 0, that can be computed using only the input wires and the output wires of the gates in layer 0, calling this layer for layer 1. We continue in this manner until all gates have been assigned a unique layer. The interesting thing to notice here is that we now

⁶ A coin-tossing protocol is needed in order to make it possible to complete a simulation proof of security, since the simulator needs to be able to extract the cut-and-choose challenge.

have a partition of the gates in such a manner that all gates in a single layer can be constructed or evaluated in parallel, in an arbitrary order, only requiring that gates at lower levels have been constructed or evaluated beforehand. Thus, given the keys of the input wires we can construct the garbled computation tables of the gates in layer 0 in an arbitrary order. Moreover, the heavy part of these computations, encryption, can be done in a SIMD manner. The only part of the construction that varies, depending on the type of gate, is which entries in the garbled computation table that should represent a 0-key and which that should represent a 1-key. Notice, however, since we implement the free XOR approach this problem is eliminated, as we can simply XOR the global key into the garbled computation table entries, which already represent a 0-key, if the output of that table entry is supposed to be a 1-key. Still, using the free XOR approach gives another problem, that is the need to further partition each layer into sets of XOR gates and non-XOR gates, in order to achieve complete SIMD *or* to keep the amount of layers and instead execute each layer like it *only* consists of XOR gates *and* execute it like it *only* consists of non-XOR gates and only use the relevant result of each of the gates. In our implementation we do the latter since preliminary tests showed this was the fastest of the two in our context.

Finally, it should be noted that the global key we choose needs to be the same for all the gates in one GC, but different for each of the GCs we make to allow opening in cut-and-choose. Keeping these changes, and this way to parallelize in mind, the protocol for construction is the same as the optimized protocol for generic GC generation previously described in Section 2, repeated ℓ' times.

The evaluation proceeds in the same manner as in the generic garbled circuit evaluation. However, we still use the same paradigm for parallelization as during construction; we evaluate each gate in a given layer, in all the evaluation circuits, in parallel, until the evaluator finds the output bits on the final layer. Finally, the evaluator takes the majority of the outputs to be his output.

Commitments. For the commitments to long strings (longer than $O(\kappa)$ bits) we can use the fact that we are in the ROM and do recursive hashing: Concatenate some random salt to the string we need to commit to and assume it is now l bits long, then we can simply hash all independent and continues $O(\kappa)$ bits pieces in parallel, resulting in $\lceil l/O(\kappa) \rceil$ digests. We concatenate these digests and again hash all independent and continues $O(\kappa)$ sized pieces of this string in parallel. We continue in this manner until a single digest remains which will be the commitment. The overhead of this approach is only logarithmic in the size of the string we hash.⁷

The OT Extension. We use the approach from [NNOB12] (see Appendix A), as the base of our OT extension. However, we make a few modifications to reduce as many operations as possible to parallel computable hashes of short bitstrings.

Assuming the existence of random oracles and a secure implementation of a κ -bit 1-out-of-2 OT as an ideal resource, the protocol is UC secure against a malicious adversary. For the rest of this section we let τ be the amount of bits in the evaluator's input for the augmented circuit, i.e. $\tau = \max(4 \cdot (|y| + |x| + \ell - 1), 8 \cdot \ell)$, where the term $|y| + |x| + \ell - 1$ comes from the first augmentation to ensure consistency of Alice's input.

Define the evaluator's (Bob's) input to the augmented circuit as a bitstring y'' of τ bits. Define $H(\cdot)$ to be a hash function with κ bits output. The modified OT extension goes as follows:

1. Bob chooses $\lceil \frac{8}{3}\kappa \rceil$ pairs of seeds, each consisting of κ random bits. That is, for each $i = 1, \dots, \lceil \frac{8}{3}\kappa \rceil$ we let $(l_i^0, l_i^1) \in_R \{0, 1\}^\kappa \times \{0, 1\}^\kappa$ be the i 'th seed pair.
2. Alice now samples $\lceil \frac{8}{3}\kappa \rceil$ random bits, $x_1, \dots, x_{\lceil \frac{8}{3}\kappa \rceil} \in_R \{0, 1\}$.
3. Alice and Bob then run $\lceil \frac{8}{3}\kappa \rceil$ OTs where, for $i = 1, \dots, \lceil \frac{8}{3}\kappa \rceil$, Bob offers (l_i^0, l_i^1) and Alice selects x_i , and receives $l_i^{x_i}$.

⁷ Notice that even if we do not assume the ROM but just that the hash function is collision free, then this block wise recursive hashing will also be a collision free hash function. This is a result due to Damgaard [Dam89].

4. Now, for each of the $i = 1, \dots, \lceil \frac{8}{3}\kappa \rceil$ pairs of random bits Bob computes the following two vectors of τ bits, using $id_{i,j}$ as a unique ID:

$$\begin{aligned} L_i^0 &= \text{H}(id_{i,0} \| l_i^0) \| \text{H}(id_{i,1} \| l_i^0) \| \dots \| \text{H}(id_{i,\tau/\kappa} \| l_i^0), \\ L_i^1 &= \text{H}(id_{i,0} \| l_i^1) \| \text{H}(id_{i,1} \| l_i^1) \| \dots \| \text{H}(id_{i,\tau/\kappa} \| l_i^1). \end{aligned}$$

5. Now, in the same manner Alice extends each of her outputs of the OT from their original length of κ bits, into strings of τ bits. Thus, Alice computes

$$L_i^{x_i} = \text{H}(id_{i,0} \| l_i^{x_i}) \| \text{H}(id_{i,1} \| l_i^{x_i}) \| \dots \| \text{H}(id_{i,\tau/\kappa} \| l_i^{x_i}).$$

6. Now, for each $i = 1, \dots, \lceil \frac{8}{3}\kappa \rceil$ Bob computes a bitstring, $\lambda_i = L_i^0 \oplus L_i^1 \oplus y''$, and sends these to Alice.
7. For each $i = 1, \dots, \lceil \frac{8}{3}\kappa \rceil$ Alice computes a bitstring as follows

$$L_i^{x_i} = L_i^{x_i} \oplus (x_i \cdot \lambda_i) = L_i^0 \oplus (x_i \cdot y'').$$

8. Alice then picks a uniformly random permutation

$$\pi : \left\{ 1, \dots, \left\lceil \frac{8}{3}\kappa \right\rceil \right\} \rightarrow \left\{ 1, \dots, \left\lceil \frac{8}{3}\kappa \right\rceil \right\}$$

where, for all i , $\pi(\pi(i)) = i$, and sends these to Bob. Furthermore, define $S(\pi) = \{i | i \leq \pi(i)\}$, that is, for each pair, the smallest index is in $S(\pi)$.

9. Now, for all the $\lfloor \frac{4}{3}\kappa \rfloor$ indexes $i \in S(\pi)$ do the following:
(a) Alice computes $d_i = x_i \oplus x_{\pi(i)}$ and sends these to Bob.
(b) Alice and Bob both compute $Z_i = \left(L_i^{x_i} \oplus L_{\pi(i)}^{x_{\pi(i)}} \right)$. This is possible for Bob since d_i uniquely determines the way to compute Z_i , i.e. if he should XOR L_i^0 with y'' .
10. For all $i \in S(\pi)$, Alice and Bob concatenate the Z_i strings, call Alice's result Z^A and Bob's result Z^B . They then check that $Z^A = Z^B$ using the following subprotocol:
(a) Alice chooses a random string $r \in_R \{0, 1\}^\kappa$.
(b) She then views her input string and r as $\left\lceil \frac{|Z^A| + \kappa}{\rho} \right\rceil$ blocks, each of ρ bits. In parallel, she then hashes each of these blocks using the hash function $\text{H}(\cdot)$.
(c) She now has $\left\lceil \frac{|Z^A| + \kappa}{\rho} \right\rceil$ hash values. She then concatenates two adjacent digests, and hash all of these.
(d) She continues in this manner, recursively concatenating and hashing the results from the previous round. At the end she has a single hash value, call it c , which she sends to Bob.
(e) Bob then sends Z^B to Alice, who checks that $Z^A = Z^B$, if this is so she then sends Z^A and r to Bob.
(f) Bob then views $Z^A \| r$ as a string of ρ bit blocks, which he, like Alice, recursively hashes and concatenates to achieve a single hash value, c' .
(g) Bob then checks if $c' = c$ and that $Z^A = Z^B$.
(h) If all checks are successful then the strings are equal, and the protocol continues, otherwise the parties abort.
11. For each $i = 1, \dots, \lfloor \frac{4}{3}\kappa \rfloor$ and for each $j = 1, \dots, \tau$ Alice defines K_j to be the string consisting of the j 'th bits from all the strings $L_i^{x_i}$, i.e.

$$K_j = L_1^{x_1}[j] \| L_2^{x_2}[j] \| \dots \| L_{\lfloor \frac{4}{3}\kappa \rfloor}^{x_{\lfloor \frac{4}{3}\kappa \rfloor}}[j].$$

This means that she gets τ keys consisting of $\lfloor \frac{4}{3}\kappa \rfloor$ bits.

12. Now, for each $i = 1, \dots, \lfloor \frac{4}{3}\kappa \rfloor$ and for each $j = 1, \dots, \tau$ Bob sets M_j to be the string consisting of the j 'th bits from all the strings L_i^0 , i.e.

$$M_j = L_1^0[j] \| L_2^0[j] \| \dots \| L_{\lfloor \frac{4}{3}\kappa \rfloor}^0[j].$$

13. Alice lets Γ_A be the string consisting of all the bits x_i for $i \in S(\pi)$, i.e. $\Gamma_A = x_1 \| x_2 \| \dots \| x_{\lfloor \frac{4}{3}\kappa \rfloor}$.
14. Bob now computes $Y_j = H(M_j)$ and achieves (Y_0, \dots, Y_τ) . He then extends each of these to ℓ' random values. That is, for each $i = 1, \dots, \ell'$ he computes $Y_j^i = H(id_{i,j} \| Y_j)$.
15. Alice computes $X_j^0 = H(K_j)$ and $X_j^1 = H(K_j \oplus \Gamma_A)$ and achieves $((X_1^0, X_1^1), \dots, (X_\tau^0, X_\tau^1))$. She then extends each of these pairs to pairs of ℓ' random values. Specifically for each $i = 1, \dots, \ell'$ she computes the following:

$$(X_j^{0,i}, X_j^{1,i}) = (H(id_{i,j} \| X_j^0), H(id_{i,j} \| X_j^1)).$$

If the parties have been honest it should be the case, that for each $i = 1, \dots, \ell'$ and $j = 1, \dots, \tau$ we have $Y_j^i = X_j^{y''[j],i}$.

Fitting It Together. After completing the modified OT extension Bob has $\tau \cdot \ell'$ key pairs of length κ . However, these keys are of course not consistent with the random keys used for the ℓ' circuits. So, for each of the $\tau \cdot \ell'$ pairs of keys Alice has, she computes the difference between the keys she achieved as a result of the modified OT extension and the actual keys to the given GCs. That, is for each $i = 1, \dots, \ell'$ and each $j = 1, \dots, \tau$ she computes

$$\begin{aligned} \delta_j^{0,i} &= X_j^{0,i} \oplus k_j^{0,i}, \\ \delta_j^{1,i} &= X_j^{1,i} \oplus k_j^{1,i} \end{aligned}$$

where $k_j^{0,i}$ is the 0-key and $k_j^{1,i}$ is the 1-key for the particular wire, j , in the particular GC, i . Alice then sends all the pairs of δ s to Bob. For each pair, Bob can only know one X value, that is, either $X_j^{0,i}$ or $X_j^{1,i}$, because of the hiding property of the OT. This means that Bob can compute exactly his choice of key, but not the other. This follows from the security of the free-XOR approach, along with the power of the random oracle for constructing $X_j^{0,i}$ and $X_j^{1,i}$, i.e. they work as one-time-pads for the keys. Thus, we get a linking between the modified OT extension and the GCs.

5 Security and Complexity

In this section we go through the complexities of our approach and sketch why our protocol is secure.

Security of the Modified OT Extension. The overall correctness and security of the modified OT extension follow from the correctness and security of the original OT extension [NNOB12] (which is secure in the random-oracle, OT-hybrid model). However, we do make several changes to this protocol. In the following we will specify these change and sketch why they do not compromise the security of the protocol. The most significant changes between the modified OT extension, and the OT extension from [NNOB12] are; the non-random construction of y'' , the use of hashing to construct the strings L_i^0 , L_i^1 and the extension of the results in Step 14 and 15.

Non-random y'' . We need y'' to be non-random in order for Bob's output of the OT extension to be consistent with his input bits for the GCs. We do this as part of the OT protocol in order to eliminate the need for sending permutation bits which would be the normal approach in order to change random OTs into OTs of specific bitstrings and choices. By embedding this in the extension itself we save some rounds of communication complexity and make the whole OT extension phase simpler by eliminating these post processing steps. Intuitively it does not compromise the security for Bob as y'' is one-time-padded with the pseudo random strings L_i^0 in Step 6 and 7. It does not compromise the security for Alice either, because of the bit switching in the end of the protocol (Step 11 and 12), along with the fact, that Alice's vector of $x_1, \dots, x_{\lfloor \frac{8}{3}\kappa \rfloor}$ is random, and thus that Bob has no idea if y'' will be XORed into a L_i^0 in Step 7.

Construction of the L_i^0 s and L_i^1 s. In [NNOB12] a pseudo random generator is used to extend a random string of length κ to a pseudo random string of length τ . However, our approach is based on invocations of hash functions on a common seed concatenated with a unique ID. This is clearly secure in the random-oracle model.

Step 14 and 15. Like for the L_i^0 s and L_i^1 s the extensions in Step 14 and 15 does not compromise security because we assume the ROM and so as long as we extend each string along with a unique ID we can assume the output is random.

Overall security assumptions. Notice that overall our protocol follows the approach for cut-and-choose based garbled circuits of [LP07]. The major theoretical difference between their (abstract) protocol and ours is the way we handle the problem of inconsistent input from the constructor. All other changes refer to practical instantiations such as the fact that we use an OT extension for all OTs, send seeds of garbled circuits instead of commitments, realize commitments using hash functions and the specific garbling scheme used.

The protocol of [LP07] is secure in the commitment-, OT-hybrid model assuming the garbling scheme is secure. For the garbling we use the scheme of [PSSW09], which was proved secure in the random oracle model. Next, notice that in [GMS08] it is shown that sending the seeds of garbled circuits will be secure when combined with the protocol of [LP07] if the hash function used is collision robust. Furthermore, we have already argued that the OT extension is secure in the random oracle, OT-hybrid model and the fact that our circuit augmentation ensures consistency of the constructor’s input or termination of the protocol independent of the evaluator’s input, except with probability at most $2^{-\ell}$. In turn we conclude that our protocol is secure in the random-oracle, OT-hybrid model.

Parallel Complexity. First see that many of the computationally heavy calculations in the protocol are hashes. Next, notice that most of these hashes are of “small” bitstrings, bounded by $O(\kappa)$, which takes place during garbling, degarbling, commitments and OT.

For garbling and degarbling we notice that the complexity becomes bounded by the length of the input to the KDF for a given gate and the depth of the circuit to securely compute. Thus, assuming access to enough parallel processors, the garbling and degarbling time will be bounded by $O(\kappa \cdot d)$ where d is the depth of the circuit to garble and under the assumption that hashing κ bits takes $O(\kappa)$ time as each gate involves hashing $O(\kappa)$ bits.

Regarding commitments the longest string we need to commit to (and thus hash) is an entire garbled circuit, which contains $3 \cdot \kappa \cdot |C|$ bits where $|C|$ denotes the amount of non-XOR gates in the circuit we wish to compute. So assuming access to enough processors the complexity is bounded by

$$O\left(\kappa \cdot \log_{O(\kappa)/\kappa}(3 \cdot \kappa \cdot |C|/\kappa)\right) = O(\kappa \cdot \log(|C|)) .$$

This is so since we need to hash $3 \cdot \kappa \cdot |C|$ bits recursively until only a single digest of κ bits remain. Thus the amount of levels in our recursive hashing is $\lceil \log_b(3 \cdot \kappa \cdot |C|/\kappa) \rceil = \lceil \log_b(3 \cdot |C|) \rceil$. The base of the logarithm, b , is determined by the amount of bits we hash sequentially and the digest size since we will reduce the bits to be hashed on the next level by a factor b after each iteration. Since we assume we hash $O(\kappa)$ bits into κ bits the base will be $b = O(\kappa)/\kappa$. Finally, we assume it takes $O(\kappa)$ time to hash $O(\kappa)$ bits sequentially, which is the time that must be spent at each level in the recursion no matter the amount of processors.

Regarding the modified OT extension notice that almost all the hashes to be computed in a given step can be done independently of each other, and thus in parallel. However, Step 10 involves hashing strings of $\max(4 \cdot \tau, 8 \cdot \ell)$ bits. Thus, assuming access to enough processors the complexity is bounded by $O\left(\kappa \cdot \log_{O(\kappa)/\kappa}(\max(4 \cdot \tau/\kappa, 8 \cdot \ell/\kappa))\right) = O(\kappa \cdot \log(\max(\tau/\kappa, \ell/\kappa)))$ by the same argumentation as in the previous paragraph when using a recursive hashing approach.

In total, assuming access to enough processors, the complexity of the whole protocol is bounded as follows:

$$O(\kappa \cdot (d + \log(|C|) + \log(\tau/\kappa)) ,$$

under the assumption that hashing $O(\kappa)$ bits takes $O(\kappa)$ time.

6 Implementation

We now describe how we constructed our implementation in CUDA in order to achieve high efficiency, based on the knowledge of the device hardware and scheduling. It should be noted that we use SHA-1 with 160 bits digest and 512 bits blocks [iosat02] as our hash function.

Gate Generation.

Kernel Structure. First, notice that we will have a case of SIMD for every circuit in ℓ' . Thus, it is obvious to have each thread in a warp processing a distinct circuit and thus having the blocks be 1-dimensional, consisting of a constant amount of warps since this structure will give us high block occupancy. Now, since preliminary tests showed that a single warp in a block achieved greater efficiency than two or more blocks in a warp we chose to have blocks consist of 32 threads. A caveat with this is that if we wish to have ℓ' not being a multiple of 32, we will need to allocate unused memory and cores and thus have SPs do useless work.

Next we notice that all gates within a single layer can be computed in arbitrary order, thus it is obvious to have one grid dimension be the amount of gates in each layer. Furthermore, as we cannot know which order the blocks will be computed in, we will need to have an iteration of kernel launches, one launch for each layer in the circuit, in order to have the output keys of the previous layer computed and ready for computing the next layer.

Regarding memory management, we first copy the seeds onto the device, and then compute the global keys for all the circuits and the 0 keys for all the input wires in all the circuits, using a unique seed for each circuit. This is done by hashing the seed along with a unique ID in order to get a “random” key (remember we assume the ROM). Afterwards, using the generated keys, we initiate a loop of kernel launches in order to compute each layer of keys and garbled computation table entries in each circuit. Between all these launches, all the currently computed keys, along with the global keys, remain in the global memory of the device so they can be used by the next kernels. Furthermore, we keep all the currently computed garbled computation tables on the device so that all the results can be copied to the host as a batch after all the kernels have finished. In order to save memory we only store the 0-key for each wire, since the 1-key can be efficient computed by simply XORing it with the appropriate global key for a given circuit.

Finally notice that the structure of the kernel for evaluation is the same as for garbling. The only difference is that before the initial launch the garbled computation table for the whole circuit is copied from the host into the global memory along with the initial input keys, one key for each of the input wires, and a description of the circuit.

Memory Coalescing. We memory coalesce all the data we use, both in the global memory and in the shared memory. As both keys and garbled computation table entries consists of 160 bits (the digest size of SHA-1), i.e. five 32-bit words, we stored all data in *segments* of $32 \cdot 5 = 160$ words. The first entry is the first word of thread 1, the second entry is the first word of thread 2, and so on up to entry 33, which then contains the second word of thread 1, entry 34 contains the second word of thread 2 and so on. Thus, all data access is coalesced in a multiple of the warp size.

The Modified OT Extension. Unlike the generation and evaluation of the GCs, the modified OT extension involves many phases, several of which are depended on the previous phases and results from interacting with the other party. This means that we cannot have a single kernel, or even a single kernel function, in order to complete all the steps of the protocol for each party.

Like we did for the GCs we coalesce all memory in blocks of 32 words. We also make segments, which consists of $5 \cdot 32 = 160$ words, such that each segment holds a coalesced hash values or a small κ bit data array, for 32 threads. For this reason we again construct kernels to use blocks of 32 threads.

Using this choice, no coalescence conversion needs to be done to use the data from the modified OT extension with our implementation of GCs. Furthermore, this choice will still keep an efficient and scalable organisation of the memory. Also, as all the data we use for computations here is completely independent, we get the possibility of only launching a single kernel for each step of the protocol in order to avoid kernel launch overhead, resulting from the iterative launching of kernels.

The kernels needed in Step 4 and 5, and Step 14 and 15, are almost the same so we only include a description of Step 4 and 5.

Steps 4 and 5. Step 4, involves hashing $2 \cdot \lceil \frac{8}{3}\kappa \rceil$ seeds τ/κ times. In order to avoid redundant data copying of L_i^0 and L_i^1 to the device when we need to construct λ_i , we compute parts of all the three vectors, L_i^0 , L_i^1 and λ_i , in each thread. That is, we include Step 6 in the kernel. To save memory usage and bandwidth we let all the 32 threads of a single block use the same pair of seeds, thus we make each thread in a block compute 160 bits of each of the three vectors L_i^0 , L_i^1 and λ_i for the same i . Next, one dimension of the grid is responsible for computing all τ bits of the three vectors, L_i^0 , L_i^1 and λ_i , and thus contains $\lceil \frac{\tau}{32 \cdot \kappa} \rceil$ threads. The other dimension of the grid is responsible for doing this for each of the $\lceil \frac{8}{3}\kappa \rceil$ vectors that need to be computed. Step 5 proceeds in the same manner, except each block only uses a single seed and each thread only computes a single digest.

Commitments Remember that commitments to long strings are done using a recursive hashing approach. In all cases the data to be hashed have been computed by the device and thus will still be in device memory in a coalesced manner. However, we also need to allocate a large result buffer in global memory that will be used to hold intermediate results. The recursive hashing procedure then consists of a loop of kernel calls where each iteration of the loop computes a level of the recursive hashes. The kernel itself is straight forward; shared memory is allocated for each segment of data to be hashed for each thread in a block (the block dimension will be 32 to match the way data is coalesced from the previous kernels). The data is then hashed and copied straight into the result buffer, still coalesced for a block of dimension 32. At the end of each iteration the result buffer and the data to be hashed switches pointers, so the old source data will be used as result buffer and the result buffer will act as the source. At the end the data in the current result buffer is copied back to the host memory.

Further Improvements. For constructing and evaluating the GCs the hash operations are clearly the main contributing time factor. However, regarding the modified OT extension it turns out, that computing the hash values on the device, barely gave an improvement in the overall execution time, and that the main contributing time factor was that of transposing bits, i.e. Steps 11 and 12. In order to achieve a significant improvement in execution time we need to implement these steps efficiently in parallel. In order to do this we need to keep the overall hardware structure and memory hierarchy in mind.

First of all, we should notice that in order to construct one word of K_j or M_j in Step 11 and 12, we need a single bit from 32 different words in $L_i^{x_i}$ or L_i^0 . In our memory organization, these are located in non-consecutive order. However, it should be noted that the remaining 31 bits of each of the 32 words are needed in the next 31 K bitstrings, K_{j+1}, \dots, K_{j+31} . Thus, depending on the caches available, it makes sense to construct the first word of $K_j, K_{j+1}, \dots, K_{j+31}$ in a batch. That is, to load 32 words of $L_i^{x_i}$ or L_i^0 and use a single bit from each of these to construct the first word of $K_j, K_{j+1}, \dots, K_{j+31}$. For this approach to be successful we need a cache of $32 \cdot 32 = 1024$ words, or 4 kilo bytes in a 32 bit system. Fortunately, this is well within the amount of shared memory on a device.

Next, consider the subprotocol for parallel equality (Step 10 of the modified OT extension). It is simple to implement on the device, by again having blocks of 32 threads and a grid of all the blocks needed to compute the individual digests. For each party we start by loading the input string into global memory and

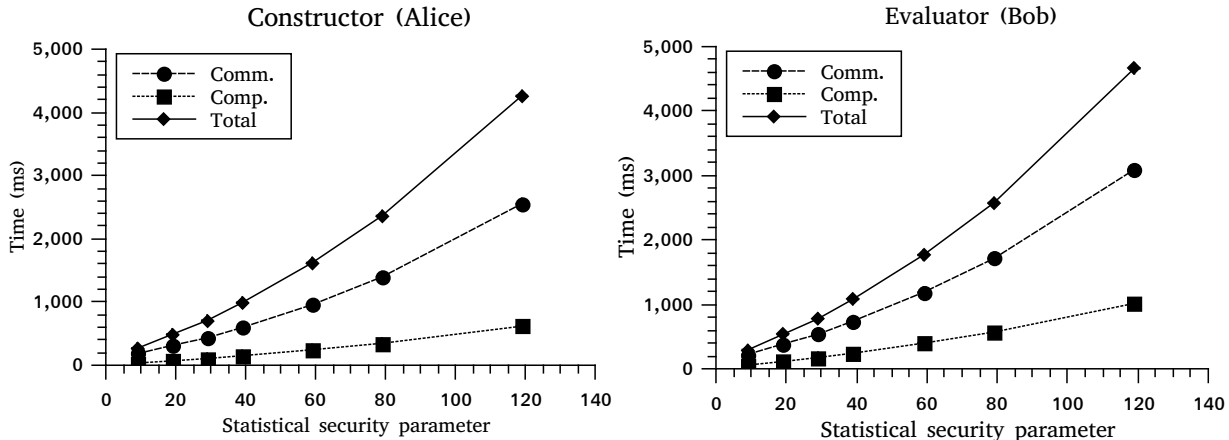


Fig. 1. Timings in milliseconds for both Alice and Bob under different statistical security parameters when computing oblivious 128 bit AES.

then construct a hash value of each, sufficiently large, chunk of bits, by loading the bits directly from global memory and storing the result back in global memory.

Using all these parallel optimizations for the modified OT extension, experiments show that the major contributing time factor is Step 15, as one might also expect, since it involved the largest amount of bits to hash.

Finally, we also introduced slight multi-threading in the host code to eliminate some idle time where one of the parties might be doing, or waiting for, network communication but still is able to do computations on the data it already has. For Alice this includes sending the difference strings in the OT extension while computing commitments to her input keys. For Bob this includes receiving and verifying consistency of Alices input keys while verifying the seeds of the check circuits.

7 Experimental Results

In this section we consider the efficiency of our protocol implementation by doing a bunch of tests. All of these tests are based on the same, commonly used, circuit for oblivious 128 bit AES encryption.⁸ This circuit is used as benchmark both in [HEKM11, LP11, NNOB12, HKS⁺10], and many more implementations of 2PC for Boolean functions. What makes this circuit a good benchmark is its relatively random structure, its relatively large size, along with its obvious practical usage, i.e. oblivious encryption.

To get the most diverse results we ran our experiments with several different statistical security parameters from 2^{-9} to 2^{-119} . We ran the experiments on two consumer grade desktop computers connected directly by a cross-over cable. At the time of purchase each of these machines had a price of less than \$1600. Both machines had similar specifications: an Intel Ivy Bridge i7 3.5 GHz quad-core processor, 8 GB DDR3 RAM, an Intel series-520 180 GB SSD drive, an MSI Z77 motherboard with gigabit LAN and an MSI GPU with an NVIDIA GTX 670 chip and 2 GB GDDR5 RAM. The machines ran the latest version (at the time) of Linux Mint. The experiments were repeated 30 times each and no front end applications were running on either of the machines and visualized in Fig. 1. These timings include every aspect of the protocol including loading circuit description and randomness along with communication between the host and device and communication between the parties. However, in the same manner as done in [NNOB12] the timing of seed OTs have not been included as this is a computation that practically only is needed once between two parties and thus will get amortized out in a practical context. The time it takes to initialize the GPU device (driver

⁸ We thank Benny Pinkas, Thomas Schneider, Nigel P. Smart and Stephen C. Williams for supplying the base circuit which we augmented for our implementation.

Table 1. Timing comparison of secure two party computation protocols evaluating oblivious 128 bit AES. d is the depth of the circuit to be computed.

	Security	ℓ	Model	Rounds	Time (s)	Equipment
[HEKM11]	Semi honest	-	ROM	$O(1)$	0.20	Desktop
This work	Malicious	2^{-9}	ROM	$O(1)$	0.29	Desktop w. GPU
This work	Malicious	2^{-29}	ROM	$O(1)$	0.78	Desktop w. GPU
[KSS12]	Malicious	2^{-80}	SM	$O(1)$	1.4	Cluster, 512 nodes
[NNOB12]	Malicious	2^{-58}	ROM	$O(d)$	1.6	Desktop
This work	Malicious	2^{-59}	ROM	$O(1)$	1.8	Desktop w. GPU
[KSS12]	Malicious	2^{-80}	SM	$O(1)$	115	Cluster, 1 node

related overhead) has *not* counted, and generally would constitute between 50 and 60 milliseconds on our test systems when the GPU is set to “persistence mode”.

The estimated total running time of the protocol is given by the total time Bob spends on the execution. This follows from the following observations:

- In the beginning of the execution, before sending the first bits, both parties do the same computation, except that Bob loads a few more bits to specify his input and Alice loads a few more bits of entropy. Thus, the maximum of the IO of the parties describe the discrepancy in their concurrent execution. It turns out, that Bob is always the one spending most time on IO.
- Alice terminates after sending the final bits of information to Bob, but Bob still needs to do more computation. Thus Bob’s time gives the upper bound on the ending of the execution.
- Finally, we don’t count the discrepancy in the wall-clock time between the time Alice and Bob start to execute the protocol program. I.e. we don’t count the time one computer waits until the other computer starts to execute the protocol.

In conclusion, Bob’s timings are good estimates for the total execution time of the protocol.

7.1 Detailed Benchmarks

Table 2 shows our detailed benchmarks, we here give some details on what exactly is counted in each cell. The “Total” timings describe the *wall-clock* timings of the execution, whereas the timings of the individual parts might be counted twice because of the multi-threaded nature of the implementation. This in particular means that the wall-clock time spent on “Circuit (comm.)” in one thread is also counted as “Circuit (comp.)” in another thread since they are computed at the same time, but on different cores.

All communication happening as part of the OT extension is counted in the row “OT (comm.)” and all other communication is counted in “Circuit (comm.)”. Everything else happening as part of the OT extension, that is not communication, is counted in “OT (comp.)”, however, only garbling, degarbling and hashing of garbled circuits are counted in “Circuits (comp.)”. All other computational aspects are reflected in the “Total” row. Furthermore, a party being idle, trying to receive or send data to the other party, but cannot do so because the other party is not yet at that step is counted in the “(comm.)” rows. This can in particular be seen in the row “OT (comm.)” and “Bob” columns.

From the timings we see that the bottleneck of the protocol is the communication complexity, even when care has been taken to do communication in parallel with computation. This becomes increasingly obvious the higher the statistical security parameter is.

8 Conclusion

We believe that our protocol approach along with the implementation yields the best practical results for malicious security two-party computation. This is so since the faster timings of [KSS12] is achieved using a

Table 2. Timing in milliseconds when computing oblivious 128 bit AES under different statistical security parameters. Uncertainty is 95% confidence intervals.

ℓ	9		19		29		39	
	Alice	Bob	Alice	Bob	Alice	Bob	Alice	Bob
IO	4.584 ± 0.03356	6.009 ± 0.2572	5.018 ± 0.31758	6.349 ± 0.3827	5.302 ± 0.2743	6.855 ± 0.48984	5.733 ± 0.44256	7.156 ± 0.5989
OT (total)	22.84 ± 6.021	41.84 ± 0.2365	19.55 ± 4.974	75.40 ± 0.9913	20.62 ± 4.962	112.5 ± 1.429	28.81 ± 6.804	156.4 ± 2.519
OT (comm.)	17.56 ± 5.784	36.96 ± 0.2135	13.52 ± 4.786	69.82 ± 1.037	13.35 ± 4.750	106.1 ± 1.468	20.09 ± 6.478	149.3 ± 2.536
OT (comp.)	5.282 ± 0.2553	4.881 ± 0.1791	6.031 ± 0.2531	5.576 ± 0.2753	7.267 ± 0.3725	6.429 ± 0.3333	8.728 ± 0.4495	7.09 ± 0.3478
Circuits (total.)	194.9 ± 5.087	239.4 ± 8.323	364.2 ± 4.069	422.0 ± 5.688	519.1 ± 3.714	598.1 ± 5.096	719.1 ± 4.661	827.6 ± 6.444
Circuits (comm.)	165.1 ± 5.072	186.3 ± 8.424	301.1 ± 4.117	315.1 ± 5.607	419.9 ± 3.625	429.7 ± 5.087	576.9 ± 4.402	592.2 ± 6.438
Circuits (comp.)	29.76 ± 0.1916	53.11 ± 0.3411	63.14 ± 1.0099	107.0 ± 0.2420	99.14 ± 1.473	168.4 ± 0.2577	142.2 ± 2.502	235.4 ± 0.1839
Total	267.4 ± 8.513	291.9 ± 8.404	489.1 ± 5.653	533.4 ± 5.668	708.5 ± 5.150	776.0 ± 5.175	990.2 ± 6.533	1082 ± 6.503

ℓ	59		79		119	
	Alice	Bob	Alice	Bob	Alice	Bob
IO	6.247 ± 0.2747	8.067 ± 0.5317	7.110 ± 0.2761	8.481 ± 0.2683	8.551 ± 0.6838	9.681 ± 0.5724
OT (total)	25.79 ± 5.286	249.7 ± 3.470	29.58 ± 5.503	351.1 ± 2.341	38.28 ± 6.408	621.1 ± 3.746
OT (comm.)	14.62 ± 5.136	240.8 ± 3.394	15.97 ± 5.583	340.6 ± 2.293	20.20 ± 6.470	606.5 ± 3.831
OT (comp.)	11.18 ± 0.5113	8.93 ± 0.2483	13.62 ± 0.5328	10.49 ± 0.3104	18.08 ± 0.5546	14.65 ± 0.2814
Circuits (total.)	1176 ± 6.736	1335 ± 6.177	1714 ± 13.261	1929 ± 6.256	3134 ± 5.685	3488 ± 8.746
Circuits (comm.)	942.3 ± 5.759	945.8 ± 5.729	1381 ± 13.108	1370 ± 5.655	2534 ± 4.253	2487 ± 8.407
Circuits (comp.)	233.8 ± 3.372	388.8 ± 1.165	333.8 ± 2.291	559.3 ± 1.291	599.6 ± 3.850	1001 ± 0.83507
Total	1619 ± 6.194	1766 ± 6.372	2363 ± 6.355	2574 ± 6.280	4262 ± 7.391	4663 ± 8.254

large grid with an estimated purchase price of at least \$129,168 per party⁹ which might not be feasible in the majority of use cases. It should further be noted that their only timings are for statistical security 2^{-80} and that we do not expect a lower security parameter to yield a significant increase in speed due to their approach in parallelization which uses one core per garbled circuit. I.e. they would not be able to utilize more than 28 or 94 cores per player if using statistical security 2^{-9} , respectively 2^{-29} . Thus using a less conservative statistical security parameter it seems highly plausible that our protocol implementation will match the pricey grid computer implementation of [KSS12].

⁹ We contacted the authors of [KSS12] who unfortunately did not have any price estimate on the Sun Blade x6240 system which they used for their timings results. Furthermore, as Sun Blade x6240 has reached end-of-life our estimate is based on the minimal price of a 256 core x86 system of the current successor of Sun Blade x6240, i.e. the Sun Blade X3-2B.

Next notice that the approach of [NNOB12] achieves a slightly faster result for a conservative statistical security parameter. However, their round complexity is asymptotically greater than ours which could yield performance issues if the protocol were to be executed on the Internet since several packet transmission must be initialized several times during the execution. Furthermore, their timings are based on amortization of 54 instances (or 27 if one is happy with statistical security 2^{-55}). Finally, by an artifact of their approach choosing a lower security parameter will not give significant performance improvements. In particular, a factor 2 in execution time seems to be the absolute maximal time improvement possible by an arbitrary reduction of the statistical security.

In conclusion, we have showed that the construction of a parallel protocol for 2PC in the SIMD PRAM model with implementation on the GPU can yield very positive results.

Future Work. Even though the time needed for the seed OTs can be amortized out from repeated use of the protocol it would still be interesting to see how fast these could be done, in particular, in parallel using the GPU.

Other interesting aspects one could try out with this protocol is to use another key derivation function, such as one based on AES. Furthermore, using our parallel protocol for covert security would also be interesting, since most of the communication complexity can be eliminated in the covert model when using seeds for generations of the garbled circuits.

Finally, implementing some working set mechanism (as done in [KSS12]) would be interesting, as it would make it possible to garble extremely large circuits.

The Code. The benchmark implementation we did for this work is freely available for non-commercial use at <http://daimi.au.dk/~jot2re/cuda>.

Acknowledgment. The authors would like to thank Roberto Trifiletti for supplying the code we used for circuit parsing and Springer for publishing the extended abstract of this work [FN13].

References

- [BCD⁺09] Peter Bogetoft, Dan Lund Christensen, Ivan Damgård, Martin Geisler, Thomas P. Jakobsen, Mikkel Krøigaard, Janus Dam Nielsen, Jesper Buus Nielsen, Kurt Nielsen, Jakob Pagter, Michael I. Schwartzbach, and Tomas Toft. Secure multiparty computation goes live. In *Financial Cryptography*, pages 325–343, 2009.
- [Bea96] Donald Beaver. Correlated pseudorandomness and the complexity of private computations. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, STOC '96, pages 479–488, New York, NY, USA, 1996. ACM.
- [Can01] R. Canetti. Universally composable security: A new paradigm for cryptographic protocols. *Foundations of Computer Science, IEEE Annual Symposium on*, 0:136, 2001.
- [Cor12] Nvidia Corporation. NVIDIA CUDA C Programming Best Practices Guide. Technical report, October 2012.
- [Dam89] Ivan Damgård. A design principle for hash functions. In *CRYPTO*, pages 416–427, 1989.
- [DNO08] Ivan Damgård, Jesper Buus Nielsen, and Claudio Orlandi. Essentially optimal universally composable oblivious transfer. In Pil Joong Lee and Jung Hee Cheon, editors, *ICISC*, volume 5461 of *Lecture Notes in Computer Science*, pages 318–335. Springer, 2008.
- [FN13] Tore Kasper Frederiksen and Jesper Buus Nielsen. Fast and maliciously secure two-party computation using the gpu. In *Applied Cryptography and Network Security - 11th International Conference, ACNS 2013, Banff, June 25-28, 2013. Proceedings*, LNCS. Springer, 2013.
- [GMS08] Vipul Goyal, Payman Mohassel, and Adam Smith. Efficient two party and multi party computation against covert adversaries. In *EUROCRYPT*, pages 289–306, 2008.
- [GMW87] Oded Goldreich, Silvio Micali, and Avi Wigderson. How to play any mental game or a completeness theorem for protocols with honest majority. In *STOC*, pages 218–229, 1987.
- [HEKM11] Yan Huang, David Evans, Jonathan Katz, and Lior Malka. Faster secure two-party computation using garbled circuits. In *USENIX Security Symposium*, 2011.

- [HKS⁺10] Wilko Henecka, Stefan Kögl, Ahmad-Reza Sadeghi, Thomas Schneider, and Immo Wehrenberg. Tasty: tool for automating secure two-party computations. In Ehab Al-Shaer, Angelos D. Keromytis, and Vitaly Shmatikov, editors, *ACM Conference on Computer and Communications Security*, pages 451–462. ACM, 2010.
- [IKNP03] Yuval Ishai, Joe Kilian, Kobbi Nissim, and Erez Petrank. Extending oblivious transfers efficiently. In *In CRYPTO 2003, Springer-Verlag (LNCS 2729)*, pages 145–161. Springer-Verlag, 2003.
- [iosat02] National institute of standards and technology. FIPS 180-2, Secure Hash Standard, Federal Information Processing Standard (FIPS), Publication 180-2. Technical report, DEPARTMENT OF COMMERCE, August 2002.
- [IPS08] Yuval Ishai, Manoj Prabhakaran, and Amit Sahai. Founding cryptography on oblivious transfer - efficiently. In *CRYPTO*, pages 572–591, 2008.
- [KH10] David B. Kirk and Wen-mei W. Hwu. *Programming Massively Parallel Processors: A Hands-on Approach (Applications of GPU Computing Series)*. Morgan Kaufmann, 1 edition, February 2010.
- [KS08] Vladimir Kolesnikov and Thomas Schneider. Improved garbled circuit: Free xor gates and applications. In Luca Aceto, Ivan Damgård, Leslie Ann Goldberg, Magnús M. Halldórsson, Anna Ingólfssdóttir, and Igor Walukiewicz, editors, *ICALP (2)*, volume 5126 of *Lecture Notes in Computer Science*, pages 486–498. Springer, 2008.
- [KSS12] Benjamin Kreuter, Abhi Shelat, and Chih-Hao Shen. Billion-gate secure computation with malicious adversaries. In *Proceedings of the 21st USENIX conference on Security symposium, Security’12*, pages 14–14, Berkeley, CA, USA, 2012. USENIX Association.
- [LOP11] Yehuda Lindell, Eli Oxman, and Benny Pinkas. The ips compiler: Optimizations, variants and concrete efficiency. In Phillip Rogaway, editor, *CRYPTO*, volume 6841 of *Lecture Notes in Computer Science*, pages 259–276. Springer, 2011.
- [LP07] Yehuda Lindell and Benny Pinkas. An efficient protocol for secure two-party computation in the presence of malicious adversaries. In *Proceedings of the 26th annual international conference on Advances in Cryptology, EUROCRYPT ’07*, pages 52–78, Berlin, Heidelberg, 2007. Springer-Verlag.
- [LP09] Yehuda Lindell and Benny Pinkas. A proof of security of yao’s protocol for two-party computation. *J. Cryptology*, 22(2):161–188, 2009.
- [LP11] Yehuda Lindell and Benny Pinkas. Secure two-party computation via cut-and-choose oblivious transfer. In Yuval Ishai, editor, *TCC*, volume 6597 of *Lecture Notes in Computer Science*, pages 329–346. Springer, 2011.
- [MNT90] Yishay Mansour, Noam Nisan, and Prasoona Tiwari. The computational complexity of universal hashing. In *Structure in Complexity Theory Conference*, page 90. IEEE Computer Society, 1990.
- [NIK12] Naoki Nishikawa, Keisuke Iwai, and Takakazu Kurokawa. High-performance symmetric block ciphers on multicore cpu and gpus. *International Journal of Networking and Computing*, 2(2), 2012.
- [NNOB12] Jesper Buus Nielsen, Peter Sebastian Nordholt, Claudio Orlandi, and Sai Sheshank Burra. A new approach to practical active-secure two-party computation. In Reihaneh Safavi-Naini and Ran Canetti, editors, *CRYPTO*, volume 7417 of *Lecture Notes in Computer Science*, pages 681–700. Springer, 2012.
- [NO09] Jesper Buus Nielsen and Claudio Orlandi. Lego for two-party secure computation. In *TCC*, pages 368–386, 2009.
- [NPS99] Moni Naor, Benny Pinkas, and Reuban Sumner. Privacy preserving auctions and mechanism design. In *ACM Conference on Electronic Commerce*, pages 129–139, 1999.
- [PDL11] Shi Pu, Pu Duan, and Jyh-Charn Liu. Fastplay-a parallelization model and implementation of smc on cuda based gpu cluster architecture. *IACR Cryptology ePrint Archive*, 2011:97, 2011.
- [PSSW09] Benny Pinkas, Thomas Schneider, Nigel Smart, and Stephen Williams. Secure two-party computation is practical. In Mitsuru Matsui, editor, *Advances in Cryptology - ASIACRYPT 2009*, volume 5912 of *Lecture Notes in Computer Science*, pages 250–267. Springer Berlin / Heidelberg, 2009. 10.1007/978-3-642-10366-7_15.
- [SS11] Abhi Shelat and Chih-Hao Shen. Two-output secure computation with malicious adversaries. In Kenneth G. Paterson, editor, *EUROCRYPT*, volume 6632 of *Lecture Notes in Computer Science*, pages 386–405. Springer, 2011.
- [XLZ11] Lei Xu, Dongdai Lin, and Jing Zou. Ecdlp on gpu. *IACR Cryptology ePrint Archive*, 2011:146, 2011.
- [Yao82] Andrew C. Yao. Protocols for secure computations. In *Proceedings of the 23rd Annual Symposium on Foundations of Computer Science, SFCS ’82*, pages 160–164, Washington, DC, USA, 1982. IEEE Computer Society.

A The OT Extension of [NNOB12]

The following protocol from [NNOB12] shows how we can get a practically unbounded amount of 1-out-of-2 κ bits OTs using only $\lceil \frac{8}{3}\kappa \rceil$ seed OTs of κ bit strings. The extension goes as follows, using ψ to denote the amount of OTs we wish to construct with P_2 denoting the player giving the input and P_1 denoting the player choosing the output:

1. P_1 samples $\Gamma_B \in_R \{0, 1\}^\psi$ and for $i = 1, \dots, \lceil \frac{8}{3}\kappa \rceil$, it also samples $L_i^0, L_i^1 \in_R \{0, 1\}^\kappa$.
2. P_2 samples $y_1, \dots, y_{\lceil \frac{8}{3}\kappa \rceil} \in_R \{0, 1\}$.
3. P_1 and P_2 then run $\lceil \frac{8}{3}\kappa \rceil$ instances of a 1-out-of-2 OT protocol in the following manner:
 - For $i = 1, \dots, \lceil \frac{8}{3}\kappa \rceil$ P_1 offers (L_i^0, L_i^1) to P_2 .
 - P_2 chooses $y_i \in \{0, 1\}$ such that it receives $L_i^{y_i}$.
 - Now, P_1 runs a pseudo random number generator on all the candidate value, (L_i^0, L_i^1) :

$$Y_i^0 = \text{prg}^\psi(L_i^0) \quad , \quad Y_i^1 = \text{prg}^\psi(L_i^1) .$$

- P_2 now runs the same pseudo random number generator on all the $\lceil \frac{8}{3}\kappa \rceil$ values of $L_i^{y_i}$ to extend each of these from their original length of κ bits into strings of ψ bits. Thus it learns $Y_i^{y_i} = \text{prg}^\psi(L_i^{y_i})$.
- Now, P_1 computes $\lceil \frac{8}{3}\kappa \rceil$ bitstrings, $\lambda_1, \dots, \lambda_{\lceil \frac{8}{3}\kappa \rceil}$, and sends these to P_2 :

$$\lambda_i = Y_i^0 \oplus Y_i^1 \oplus \Gamma_B .$$

- For $i = 1, \dots, \lceil \frac{8}{3}\kappa \rceil$ P_2 computes a semi authenticated bit as follows

$$|y_i\rangle = (y_i, Y_i^{y_i} \oplus (y_i \cdot \lambda_i)) = (y_i, Y_i^0 \oplus (y_i \cdot \Gamma_B))$$

4. P_2 now picks a uniformly random permutation:

$$\pi : \left\{ 1, \dots, \left\lceil \frac{8}{3}\kappa \right\rceil \right\} \rightarrow \left\{ 1, \dots, \left\lceil \frac{8}{3}\kappa \right\rceil \right\} ,$$

such that for all i ; $\pi(\pi(i)) = i$. That is, the permutation is a pairing. P_2 sends this pairing to P_1 . Given this pairing, let $\mathcal{S}(\pi) = \{i | i < \pi(i)\}$. That is, $\mathcal{S}(\pi)$ contains the $\lfloor \frac{4}{3}\kappa \rfloor$ elements which have the smallest index of all pairs in π .

5. Now, for all the $\lfloor \frac{4}{3}\kappa \rfloor$ indices, $i \in \mathcal{S}(\pi)$, do the following:
 - P_2 announces $d_i = y_i \oplus y_{\pi(i)}$.
 - P_1 and P_2 computes

$$|z_i\rangle = |y_i\rangle \oplus |y_{\pi(i)}\rangle \oplus d_i = (z_i, Z_i) .$$

- Notice that P_1 can simply compute this as:

$$|z_i\rangle = (y_i \oplus y_{\pi(i)} \oplus d_i, Y_i^0 \oplus Y_{\pi(i)}^0 \oplus (d_i \cdot \Gamma_B)) = (z_i, Z_i) .$$

- P_1 and P_2 concatenate their bitstrings Z_i for all $i \in \mathcal{S}(\pi)$. Call P_1 's concatenation Z^1 and P_2 's concatenation Z^2 . The parties then check equality of these strings using the following subprotocol:

- P_1 chooses a random string $r \in_R \{0, 1\}^\kappa$, computes $c = H(Z^1 || r)$ and sends c to P_2 .
- P_2 sends Z^2 to P_1 .
- P_1 check that $Z^1 = Z^2$ and if so sends Z^1 and r to P_2 , otherwise it aborts.
- P_2 checks that $H(Z^1 || r) = c$ and that $Z^1 = Z^2$, if not, it aborts.
- If no party aborts then $Z^1 = Z^2$.

6. P_1 then defines x_j be the j 'th bit of Γ_B and M_j to be the string consisting of the j 'th bits from all the strings Y_i^0 , i.e. $M_j = Y_j^0[1] || Y_j^0[2] || \dots || (Y_j^0[\lfloor \frac{4}{3}\kappa \rfloor])$. So, because of the use of the pseudo random number generator we get ψ bits, x_j , and ψ bitstrings, M_j .

7. P_2 then lets Γ_A be the string consisting of all the bits, y_i , i.e. $\Gamma_A = y_1 \parallel y_2 \parallel \dots \parallel y_{\lfloor \frac{4}{3}\kappa \rfloor}$ and lets K_j be the string consisting of the j 'th bits from all the strings $Y_i'^0 \oplus (y_i \cdot \Gamma_B)$, i.e. $K_j = (Y_1'^0 \oplus (y_1 \cdot \Gamma_B))[j] \parallel (Y_2'^0 \oplus (y_2 \cdot \Gamma_B))[j] \parallel \dots \parallel (Y_{\lfloor \frac{4}{3}\kappa \rfloor}'^0 \oplus (y_{\lfloor \frac{4}{3}\kappa \rfloor} \cdot \Gamma_B))[j]$.

8. P_1 now uses the hash function:

$$Y_j = H(M_j)$$

on each of her ψ bitstrings. Thus, he ends up having ψ strings of κ bits, (Y_1, \dots, Y_ψ) along with Γ_B .

9. P_2 also uses the hash function:

$$X_{i,0} = H(K_i) \quad , \quad X_{i,1} = H(K_i \oplus \Gamma_A),$$

twice on each of her ψ bitstrings, and ends up with ψ pairs of bitstrings, $((X_1^0, X_1^1) \dots, (X_\psi^0, X_\psi^1))$.

Now it will be the case for $i = 1, \dots, \psi$ that if $x_i = 0$ then $Y_i = X_i^0$ and if $x_i = 1$ then $Y_i = X_i^1$, i.e. a random OT.