

From oblivious AES to efficient and secure database join in the multiparty setting^{*}

Sven Laur^{2,3}, Riivo Talviste^{1,2}, and Jan Willemsen^{1,3}

¹ Cybernetica, Ülikooli 2, Tartu, Estonia

² Institute of Computer Science, University of Tartu, Liivi 2, Tartu, Estonia

³ Software Technology and Applications Competence Center, Ülikooli 2, Tartu, Estonia

Abstract. AES block cipher is an important cryptographic primitive with many applications. In this work, we describe how to efficiently implement the AES-128 block cipher in the multiparty setting where the key and the plaintext are both in a secret-shared form. In particular, we study several approaches for AES S-box substitution based on oblivious table lookup and circuit evaluation. Given this secure AES implementation, we build a universally composable database join operation for secret shared tables. The resulting protocol scales almost linearly with the database size and can join medium sized databases with 100,000 rows in few minutes, which makes many privacy-preserving data mining algorithms feasible in practice. All the practical implementations and performance measurements are done on the SHAREMIND secure multiparty computation platform.

1 Introduction

Many information systems need to store and process private data. Encryption is one of the best ways to assure confidentiality, as it is impossible to learn anything from encrypted data without knowledge of the private key. However, the number of processing steps one can carry out on encrypted data is rather limited unless we use fully homomorphic encryption. Unfortunately, such encryption schemes are far from being practical even for moderate-sized data sets [22].

Another compelling alternative is share-computing, since it assures data confidentiality and provides a way to compute on secret shared data, which is several magnitudes more efficient than fully homomorphic encryption. In this setting, data is securely shared among several parties so that individual parties learn

^{*} This research was supported by the ERDF through EXCS and STACC; the ESF Doctoral Studies and Internationalisation Programme DoRa and by Estonian institutional research grant IUT2-1.

This research was, in part, funded by the U.S. Government. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Government. Distribution Statement A (Approved for Public Release, Distribution Unlimited).

nothing about shared values during the computations and the final publication of output shares reveals only the desired output(s). For most share-computing systems, even a coalition of parties cannot learn anything about private data unless the size of a coalition is over a threshold.

Development and implementation of such multi-party computing platforms is an active research area. FairPlayMP [6], SecureSCM [3], SEPIA [14], SHAREMIND [9], VMCrypt [31] and TASTY [25] computing platforms represent only some of the most efficient implementations and share-computing has been successfully applied to real-world settings [11,10].

Note that various database operations are particularly important in privacy-preserving data processing. Efficient and secure protocols for most key operations on secret-shared databases are already known, see [30]. The most notable operation still missing is database join based on secret-shared key columns. This operation can be used e.g. for combining customer data coming from different organisations or linking the results of statistical polls into a single dataset.

Our main theoretical contribution is an efficient multi-party protocol for database join, which combines oblivious shuffle with pseudorandom function evaluation on secret-shared data. In practice, we instantiate the pseudorandom function with the AES-128 block cipher and implement it on the SHAREMIND platform [9]. The latter is a non-trivial task, since the input and the secret key are secret-shared in this context. The resulting AES-evaluation protocol is interesting in its own right. First, AES is becoming a standard performance benchmark for share-computing platforms [19,26,34,29] and thus we can directly compare how well the implementation on the SHAREMIND platform does. Second, a secret-shared version of AES can be used to reduce security requirements put onto the key management of symmetric encryption [19]. In brief, we can emulate trusted hardware encryption in the cloud by sharing a secret key among several servers.

2 Preliminaries

AES. Advanced Encryption Standard (AES) is a symmetric block cipher approved by the National Institute of Standards and Technology [32]. AES takes a 128-bit block of plaintext and outputs 128 bits of corresponding ciphertext. AES can use cipher keys with lengths of 128, 192 or 256 bits. In our work we will only use AES-128, which denotes AES with 128-bit keys.

Sharemind platform. SHAREMIND platform is a practical and secure share-computing framework for privacy-preserving computations [9], where the private data is shared among three parties referred to as *miners*. In its original implementation, SHAREMIND uses additive secret sharing on 32-bit integers, i.e., a secret s is split into three shares s_1, s_2, s_3 such that $s = s_1 + s_2 + s_3 \bmod 2^{32}$. In this work, we use bitwise sharing where the secret can be reconstructed by XOR-ing individual shares: $s = s_1 \oplus s_2 \oplus s_3$.

The current SHAREMIND implementation is guaranteed to be secure only if the adversary can observe the internal state of a single miner node. Thus, we report performance results only for the *semi-honest setting*. Additionally, we show how to generalise our approach to malicious setting. The latter is rather straightforward, as all protocols are based only on secure addition and multiplication protocols. Although the bitwise sharing alone is not secure against *malicious corruption*, shared message authentication codes can be used to guarantee integrity of secret sharings throughout the computations [20,33].

Security definitions and proofs. We use standard security definitions based on ideal versus real world paradigm. In brief, security is defined by comparing a real protocol with an ideal implementation where a trusted third party privately collects all inputs, does all computations and distributes outputs to corresponding parties. We say that a protocol is secure if any plausible attack against real protocol can be converted to an attack against ideal protocol such that both attacks have comparable resource consumption and roughly the same success rate, see standard treatments [23,15,16] for further details.

A canonical security proof uses a wrapper (*simulator*) to link a real world adversary with the ideal world execution model. More precisely, the simulator has to correctly fake missing protocol messages and communicate with the trusted party. As most protocols are modularly built from sub-protocols, security proofs can be further compacted. Namely, if all sub-protocols are *universally composable*, then we can prove the security in the hybrid model where executions of all sub-protocols are replaced with ideal implementations [16].

Since almost all share-computing platforms including SHAREMIND provide universally composable data manipulation operations, we use this compossibility theorem to omit unnecessary details from security proofs (see also [9]).

Efficiency metrics in protocol design. Real-life efficiency of a protocol execution depends on the number of rounds and the total amount of messages sent over communication channels. The actual dependency is too complicated to analyse directly. Hence, we consider two important sub-cases. When the total communication is small compared to channel bandwidths, then the running time depends linearly on the number of rounds. If the opposite holds, then running time depends linearly on the communication complexity.

3 Share-computing protocol for AES block cipher

The overall structure of our protocol follows the standard AES algorithm specification [32]. However, there are some important differences stemming from the fact that the secret key and the message is bitwise secret shared and we have to use share-computing techniques. Fortunately, three out of four sub-operations are linear and thus can be implemented by doing local share manipulations. The efficiency of the AES protocol implementation is determined by `SubWord()`

and `SubBytes()` operations that evaluate the S-box on secret-shared data. The `SubWord()` function used in key expansion applies the S-box independently to each byte of its input word. Similarly, the `SubBytes()` function uses the S-box independently on each byte of the 4-word state given as the argument.

3.1 S-box evaluation protocol based on oblivious selection

As the AES S-box is a non-linear one-to-one mapping of byte values, it can be implemented as 256 element lookup table. In our setting, the input of the S-box is secret shared and we need oblivious array selection to get the shares of the right table entry. The latter can be achieved by using various techniques from [30]. First, we must convert the input x into a zero-one index vector \mathbf{z} where all entries, except one, are zeros. The non-zero vector element z_x corresponds to the entry in the S-box array that we want to pick as the output. More precisely, let $x_7x_6 \dots x_0$ be the bit-representation of the input x and $i_7i_6 \dots i_0$ be the bit-representation of an index i . Then $z_i = [x_7 = i_7] \wedge \dots \wedge [x_0 = i_0]$ and the shares of index vector \mathbf{z} can be computed by evaluating multinomials

$$z_i = (x_7 \oplus i_7 \oplus 1) \cdots (x_0 \oplus i_0 \oplus 1) . \quad (1)$$

For example, the first entry can be computed as $z_0 = (1 \oplus x_7)(1 \oplus x_6) \dots (1 \oplus x_0)$ and the second entry as $z_1 = (1 \oplus x_7)(1 \oplus x_6) \dots (1 \oplus x_1)x_0$.

Note that each multinomial z_i is of degree 8 and thus 1792 secure multiplications over \mathbb{F}_2 are needed. To reduce the number of communication rounds, we gather terms $b_{ij} = x_j \oplus i_j \oplus 1$ into eight 256 element vectors:

$$\mathbf{b}_7 = (b_{0,7}, \dots, b_{255,7}), \dots, \mathbf{b}_0 = (b_{0,0}, \dots, b_{255,0})$$

and use vectorised bitwise multiplications to multiply all eight terms in the same row. If we do them sequentially, then the computation of index vector requires seven multiplication rounds. With tree-style evaluation strategy we can reduce the number of multiplication rounds to three. For that, we must evaluate same level brackets in parallel for $\mathbf{z} = ((\mathbf{b}_7 \cdot \mathbf{b}_6) \cdot (\mathbf{b}_5 \cdot \mathbf{b}_4)) \cdot ((\mathbf{b}_3 \cdot \mathbf{b}_2) \cdot (\mathbf{b}_1 \cdot \mathbf{b}_0))$. The multiplicative complexity of this step can be further decreased by utilising the underlying recursive structure of the index vector, as proposed by Launchbury *et al.* [29]. For comparison, we also reimplemented their solution.

As the second step, we must compute scalar product between the indicator vector \mathbf{z} and 256-element output table \mathbf{y} of the S-box. As elements of \mathbf{y} are 8-bit long whereas elements of \mathbf{z} are from \mathbb{F}_2 , we must select output bits one by one. Let $\mathbf{y}_j = (y_{0,j}, \dots, y_{255,j})$ denote the vector of j th bits in the output table \mathbf{y} . Then the j th output bit f_j of the S-box can be computed as

$$f_j = \langle \mathbf{z}, \mathbf{y}_j \rangle = \sum_{i=0}^{255} z_i y_{ij} \quad (2)$$

over \mathbb{F}_2 . Since the output table \mathbf{y} is public, all operations can be done locally and the second step does not contribute to the communication complexity.

3.2 S-box evaluation protocol based on circuit evaluation

The oblivious indexing as a generic approach is bound to provide a protocol with sub-optimal multiplication complexity, as the two stage evaluation of output bits f_j forces us to compute terms z_i that are dropped in the equation (2).

We can address this issue by secure computation techniques based on branching programs [27]. For that, we must convert the expression for f_j into a binary decision diagram \mathcal{B} with minimal number of decision nodes. After that we must build a corresponding arithmetic circuit that evaluates \mathcal{B} in bottom-up manner. As each decision node introduces two secure multiplications, the efficiency of the resulting protocol is determined by the shape of \mathcal{B} . Let c denote the total number of decision nodes and d denote the longest path in \mathcal{B} . Then the resulting protocol consists of $2c$ secure multiplication operations over \mathbb{F}_2 , which can be arranged into d rounds of parallel multiplications.

Although this approach produces significant gains, we can use recent findings in hardware optimisation to boost efficiency further. Circuit minimisation for the AES S-box is a widely studied problem in the hardware design with many known results. In this work, we use the designs by Boyar and Peralta [12,13]. Note that their aim was to minimise the total number of gates and the overall circuit depth, while we need a circuit with minimal number of multiplication gates (AND operations) and with paths that contain as few multiplications as possible, i.e., have a low multiplicative depth. Hence, their best design with 128 gates is not the best for our purposes, as it contains 34 multiplications and its multiplicative depth is 4, while their older design [12] contains 32 multiplication and has a multiplicative circuit depth 6. Of course, the multiplicative depth plays also important role in the protocol, when the bandwidth is high, hence, the newer design might have advantages when only a few AES evaluations are performed.

As extended versions of both articles contain straight-line C-like programs for their circuits, it is straightforward to implement the corresponding secure evaluation protocol with a minor technical tweak. As byte is the smallest data unit supported by network communication libraries, entire byte is used to send elements of \mathbb{F}_2 over the network during a secure multiplication protocol. We can eliminate this bloat by doing eight multiplications in parallel, since eight individual values can be packed into the same byte.

It is straightforward to achieve this grouping for the `SubBytes()` function, as it evaluates 16 S-boxes in parallel. Consequently, if we treat original variables as 16-element bit-vectors, we can evaluate all 16 copies of the original circuit in parallel without altering the straight-line program. For the `SubWord()` function, additional regrouping is necessary, as it evaluates only four S-boxes in parallel. It is sufficient if we must split all multiplications into pairs that can be executed simultaneously so that we can do eight multiplications in parallel.

3.3 Security analysis for the entire protocol

Note that all three versions of the AES S-box evaluation algorithms are arithmetic circuits consisting of addition and multiplication gates. Hence, it is straightforward to prove the following result.

Theorem 1. *If a share-computing framework provides universally composable protocols for bitwise addition, bitwise multiplication and bit decomposition, then all three AES S-box implementations are universally composable. Any universally composable AES S-box implementation gives a rise to a universally composable share-computing protocol for the AES block cipher.*

Proof. The proof follows directly from the universal composability theorem as we use share-computing protocols to evaluate arithmetic circuits. \square

Note that this result holds for any corruption model including the SHAREMIND framework, which provides security against one-out-three static passive corruption. To get security against active corruption, the underlying secret sharing scheme must support both bitwise addition and multiplication while being verifiable. There are two principal ways to achieve this.

First, we can embed elements of \mathbb{F}_2 into some larger finite field \mathbb{F}_{2^t} with extension element α and then use standard verifiable secret sharing schemes which support secure multiplication over \mathbb{F}_{2^t} . On top of that it is rather straightforward to implement universally composable bit decomposition [18], which splits a secret $x \in \mathbb{F}_{2^t}$ into a vector of shared secrets x_{t-1}, \dots, x_0 such that $x = x_{t-1}\alpha^{t-1} + \dots + x_1\alpha + x_0$. As a consequence, all three assumptions of Theorem 1 are satisfied and we get a secure protocol for evaluating AES. However, there is a significant slowdown in the communication due to prolonged shares.

Alternatively, we can use oblivious message authentication [20] to protect individual bits without extending shares. However, this step attaches a long secret shared authentication code to each bit. To avoid slowdown, we can authenticate long bit vectors with a single authentication code. The latter fits nicely into the picture, as we have to evaluate 16 circuits in parallel.

3.4 Further tweaks of the AES evaluation protocol

Block ciphers are often used to encrypt many messages under the same secret key. In such settings, it is advantageous to encrypt several messages in parallel in order to reduce the number of communication rounds. The latter is straightforward in the SHAREMIND platform, as it naturally supports parallel operations with vectors. The corresponding vectorised AES protocol takes in a vector of plaintext shares and a vector of shared keys and outputs a vector of cipher text shares. As another efficiency tweak note that we need to execute that key scheduling only once if the secret key is fixed during the encryption. Hence, we can run the key scheduling protocol separately and store the resulting shares of all 128-bit round keys for later use. The corresponding separation of pre-processing and online phases decreases amortised complexity by a fair margin.

Protocol	Multiplicative depth	Running time (1 evaluation)	Multiplicative complexity	Running time (4096 evaluations)
OBSEL	3	32.5 ms	1792	9051 ms
LDDAM	3	31.1 ms	304	1109 ms
BCIRC-1	6	69.6 ms	32	148 ms
BCIRC-2	4	40.8 ms	34	127 ms

Table 1. Performance results of various S-box evaluation algorithms.

3.5 Efficiency metrics and real-life performance

Having established essentially four methods with very different complexity parameters, we need to compare their real-life performance. For that we have implemented four versions of `SubBytes()` routines on the SHAREMIND platform and measured the actual performance. The tests were done on a cluster where each of the three SHAREMIND miners was deployed in a separate machine. The computers in the cluster were connected by an ethernet local area network with link speed of 1 Gbps. Each computer in the cluster had 48 GB of RAM and a 12-core 3 GHz CPU with Hyper Threading. The channels between the computers were also encrypted using 256-bit elliptic curve key agreement and the ChaCha stream cipher [8] provided by the underlying RakNet networking library [2]. While the choice of ChaCha is not standard, the best known attacks against it are still infeasible in practice [5].

We considered algorithms in two different settings. First, we measured the time needed to complete a single evaluation of `SubBytes()` function. Second, we measured how much time does it take to evaluate 4096 `SubBytes()` calls in parallel. The first setting corresponds to the case where various delays have dominant impact on the running-time, whereas the effect of communication complexity dominates in the second case. Table 1 compares theoretical indicators⁴ and practical performance for all four protocols. The OBSEL protocol is based on oblivious selection vector and LDDAM is the same protocol with reduced number of multiplications [29]. Protocols based on Boolean circuits designed by Boyar and Peralta are denoted by BCIRC-1, BCIRC-2.

The results clearly show that multiplicative depth and complexity are good theoretical performance measures for optimising the structure of arithmetic circuits, as they allow us to predict the running times with 10 – 20% precision. Each communication round costs 10 – 12 ms in single operation mode and each multiplication operation adds 3.5 – 5.1 ms to amortised running-time.

Secondly, we measured amortised cost of the AES evaluation protocol with precomputed round keys, see Figure 1. As expected, various algorithms have different saturation points where further parallelisation does not decrease the amortised cost any more. In particular, note that for few blocks the amortised costs of LDDAM and circuit evaluation algorithms BCIRC-1 and BCIRC-2 is

⁴ As all multiplications are carried over \mathbb{F}_2 , we do not have to compensate for various input lengths and can just count the number of multiplications.

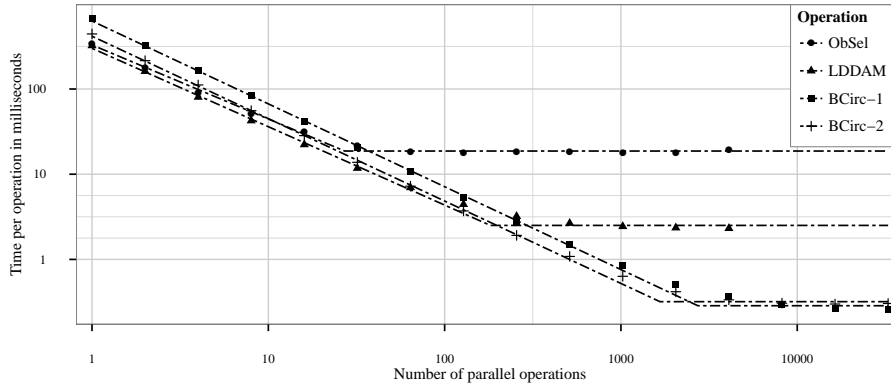


Fig. 1. Performance of AES evaluation protocols using precomputed round keys.

comparable, i.e., the advantage of circuit evaluation manifests only if we encrypt around 80 plaintexts in parallel. Also, note that the newer design BCIRC-2 with smaller multiplicative depth performs better when the number of encryption calls is between 100 – 10,000. After that the impact of communication complexity becomes more prevalent and the BCIRC-1 protocols becomes more efficient.

As the final test, we measured the running time of the AES protocol with and without key scheduling. Table 2 depicts the corresponding results. As before, we give the running times for a single encryption operation and limiting cost of a single operation if many encryptions are done in parallel. Mode I denotes encryption with key expansion and mode II denotes encryption with pre-expanded secret key. Again, the results are in good correspondence. The cost of a single operation is roughly two times slower with the key expansion⁵, since computing a shared round key requires one parallel invocation of S-boxes. For the amortised cost, the theoretical speedup should be 1.25 as there are 20 S-box invocation per round in the mode I and 16 invocations per round in the mode II. The difference in actual speedup factors suggest existence of some additional bottlenecks in our key-expansion algorithms.

Table 3 compares our results with the state of the art in oblivious AES-128 evaluation protocols. To make results comparable, the table contains results only for the semi-honest setting. In most cases, authors report the performance of AES with pre-shared keys (mode II). More than tenfold difference between two-party and three-party implementations is expected, as two-party computations require costly asymmetric primitives. Note that the cost of single operation for our implementation in Table 3 uses the approach of Launchbury *et al.*, whereas the amortized time is obtained using the circuit-based approach.

⁵ The slowdown can be further reduced to 1.2 if we compute next subkey in parallel with the AES round to reduce multiplicative depth of the circuit.

	Single operation			Amortised cost		
	Mode I	Mode II	Ratio	Mode I	Mode II	Ratio
OBSEL	682 ms	343 ms	1.99	20.34 ms	18.69 ms	1.09
LDDAM	652 ms	323 ms	2.02	4.16 ms	2.51 ms	1.66
BCIRC-1	1329 ms	664 ms	2.00	0.48 ms	0.29 ms	1.68
BCIRC-2	890 ms	443 ms	2.01	0.37 ms	0.32 ms	1.17

Table 2. Performace results for various AES evaluation algorithms

We can not fully explain roughly 20 times performance difference between the two implementations of single operation following the approach of Launchbury *et al.* Possible explanations include measurement error and extreme concentration on the network layer optimization by the authors of [29].

Authors	Reference	Setting	Mode	Single operation	Amortised cost
Pinkas et al.	[34]	2-party	II	5000 ms	— ms
Huang et al.	[26]	2-party	II	200 ms	— ms
Damgård and Keller	[19]	3-party	I	2000 ms	— ms
Launchbury et al.	[29]	3-party	II	14.28 ms	3.10 ms
This work		3-party	II	323 ms	0.29 ms

Table 3. Comparison of various secure AES-128 implementations

4 Secure database join

As mentioned in the introduction, secure database join is a way to combine several data sources in privacy-preserving manner. In this work, we consider the most commonly used *equi-join*⁶ operation, which merges tables according to one of few key columns using the equality comparison in the join predicate. In many cases, the key value is unique, such as social security number or name and postal code combined. The uniqueness assumption significantly simplifies our task. The need to deal with the colliding keys significantly increases the complexity of the protocols, and this case is handled in the extended version of the paper [1].

An ideal secure inner join protocol takes two or more secret-shared database tables and produces a new randomly ordered secret-shared table that contains the combined rows where the join predicate holds. The parties should learn nothing except for the number of rows in the new database. The random reordering of the output table is necessary to avoid unexpected information propagation when some entries are published either for input or for the output table.

⁶ The authors adapt the Structured Query Language (SQL) terminology in this paper.

Database shuffling phase

1. Miners obliviously shuffle each database table T_i .
Let T_i^* denote the resulting shuffled table with a key column \mathbf{k}_i^* .

Encryption and join phase

2. Miners choose a pseudorandom permutation π_s by generating a shared key s .
3. Miners obliviously evaluate π_s on all shared key columns \mathbf{k}_i^* .
4. Miners publish all values $\pi_s(k_{ij}^*)$ and use standard database join to merge the tables based on columns $\pi_s(\mathbf{k}_i^*)$. Let T^* be the resulting table.

Optional post-processing phase for colliding keys

5. If there are some non-unique keys in some key column $\pi_s(\mathbf{k}_i^*)$, miners should perform additional oblivious shuffle on the secret-shared table T^*

Protocol 1: Secure implementation of PRPJOIN operation

Let m_1 and m_2 denote the number of rows and n_1 and n_2 the number of columns in the input tables. Then it is straightforward to come up with a solution that uses $\Theta(m_1 m_2)$ oblivious comparison operations by mimicking a naïve database join algorithm. We can obliviously compare all the possible key pairs, shuffle the whole database, open the comparison column and remove all the rows with the equality bit set to 0. It is straightforward to prove that this protocol is secure, since it mimics the actions of ideal implementation in verbatim. We will refer to this algorithm as NAIVEJOIN and treat it as a baseline solution.

4.1 Secure inner join based on unique key column

As the first step towards a more efficient algorithm, consider a setting where the computing parties (miners) obliviously apply pseudorandom permutation π_s to encrypt the key column. As π_s is a pseudorandom permutation (a block cipher depending on an unknown key s) and all the values in the key column are unique, the resulting values look completely random if none of the miners knows π_s . Hence, it is secure to publish all the encryptions of key columns. Moreover, the tables can be correctly joined using the encryptions instead of key values.

However, such a join still leaks some information – miners learn which database rows in the first table correspond to the database rows in the second table. By shuffling the rows of initial tables, this linking information is destroyed. The resulting algorithm is depicted as Protocol 1. We emphasise that in each step all the tables are in secret-shared form. In particular, each miner carries out step 4 with its local shares and thus the table T^* is created in a secret-shared form. Note that we have also added step 5 to deal with the case of colliding keys. This case will be discussed in the extended version of the paper [1].

As the actual join operation is performed on public (encrypted) values, the construction works also for the *left* and *right outer joins*, where either the left or right table retains all its rows, whether a row with a matching key exists in the other table or not. These outer joins are common in data analysis. For instance,

given access to supermarket purchases and demographic data, we can use outer join to add person’s wealth and his/her home region to each transaction, given that both tables contain social security number. As the data about some persons might be missing from the demographic database, miners must agree on predefined constants to use instead of real shares if the encrypted key is missing. In this case, optional post-processing step is needed to hide rows with dummy values. However, the post-processing phase does not hide the number of missing data entries. We discuss this issue in the extended version of the paper [1].

Theorem 2. *Let $\mathcal{P} = (\pi_s)$ be a pseudorandom permutation family. If a share-computing framework provides universally composable protocols for database shuffle and oblivious evaluation of $\pi_s(x)$ from secret shared values of x and s , and there are no duplicate key values in any of the input tables, then the PRPJOIN protocol is universally composable in the computational model.*

Proof (Sketch). For clarity, let us analyse the security in the modified setting where \mathcal{P} is the set of all permutations and Steps 1–4 are performed by trusted third party. Let m be the number of rows in the final database table and \mathbf{y}_1 and \mathbf{y}_2 the vectors of encrypted values published during PRPJOIN protocol. For obvious reasons, $|\mathbf{y}_1 \cap \mathbf{y}_2| = m$ and the set $\mathbf{y}_1 \cup \mathbf{y}_2$ consists of $m_1 + m_2 - m$ values, which are chosen randomly from the input domain without replacement. As Step 1 guarantees that the elements in \mathbf{y}_1 and \mathbf{y}_2 are in random order, it is straightforward to simulate \mathbf{y}_1 and \mathbf{y}_2 given only the number of rows m .

Hence, the simulation of the protocol is straightforward. First, the simulator forwards all input shares and gets back the final output shares and thus learns m . After that it generates shares for the shuffled databases by creating the correct number of valid shares of zero. As the adversarial coalition is small enough, the adversary cannot distinguish them from valid shares. Next, it generates \mathbf{y}_1 and \mathbf{y}_2 according to the specification given above and forwards the values to the adversary together with properly aligned output shares such that a semihonest adversary would assemble the database of output shares in the correct way.

It is easy to see that the simulation is perfect in the semihonest model. The same is true for the malicious model with honest majority, since honest parties can always carry out all the computations without the help from the adversarial coalition. In case of dishonest majority, the adversarial coalition is allowed to learn its output and then terminate the protocol. In our case, the simulator must terminate the execution when the adversarial coalition decides to stop after learning the encrypted vectors \mathbf{y}_1 and \mathbf{y}_2 .

We can use the same simulation strategy for the original protocol where the trusted third party uses a pseudorandom permutation family. As the key s is unknown to all parties, the joint output distributions of the real and hybrid worlds are computationally indistinguishable. The latter is sufficient, as security in the hybrid model carries over to the real world through universal composable of share shuffling and oblivious function evaluation protocols. \square

Efficiency. By combining the secure oblivious AES evaluation and the oblivious shuffle from [30], we get an efficient instantiation of the PRPJOIN protocol. For

Offline phase

1. Generate shared random keys (k_{ij}) for the Carter-Wegman construction.

Online hashing phase

2. Treat each key tuple as a long bit string $\mathbf{x} = (x_s, \dots, x_1)$.
3. Use secure scalar product algorithm to compute the secret shared hash code:

$$h(\mathbf{k}_j, \mathbf{x}) = x_s k_{sj} + \dots + x_1 k_{1j} .$$

Protocol 2: Oblivious hashing OHASH

all database sizes, the resulting protocol does $\Theta(m_1 + m_2)$ share-computing operations and $\Theta(m_1 \log m_1 + m_2 \log m_2)$ public computation operations.⁷

4.2 Secure inner join based on unique multi-column key values

Let us now consider the case when database tables are joined based on several columns, such as name and birth date. We can reduce this kind of secure join to the previous case by using oblivious hashing. An ε -almost universal hash function is a function $h : \mathcal{K} \times \mathcal{M} \rightarrow \mathcal{T}$ that compresses message into shorter tags so that the following inequality holds:

$$\forall x \neq x' \in \mathcal{M} : \Pr [k \leftarrow \mathcal{K} : h(k, x) = h(k, x')] \leq \varepsilon.$$

Such a function can be used to reduce the length of the unique key that spans over several columns. However, this function must support efficient oblivious evaluation. The Carter-Wegman construction [17]

$$h(\mathbf{k}, \mathbf{x}) = x_s k_s + \dots + x_2 k_2 + x_1 k_1$$

is a good candidate for our application as it consist of a few simple operations and it is $2^{-\ell}$ almost universal when computations are done over the field \mathbb{F}_{2^ℓ} . Another compelling alternative is to use several independent Carter-Wegman functions over \mathbb{F}_2 . For ℓ independently chosen keys, the collision probability is still $2^{-\ell}$. In the semihonest model, the communication complexity of the resulting oblivious hashing protocols is the same, as the amount of communication scales linearly wrt the bit length. For the malicious models, the trade-offs depend on exact implementation details of multiplication protocol. The resulting algorithm for oblivious hashing is depicted as Protocol 2.

Theorem 3. *If a share-computing framework provides universally composable protocols for addition and multiplication over \mathbb{F}_2 , the OHASH protocol is universally composable in the information theoretical model. For ε -almost universal*

⁷ The theoretical asymptotic complexity is higher, as the size of the database can be only polynomial in the security parameter and thus oblivious PRF evaluation takes $\text{poly}(m)$ steps. Consequently, the protocol is asymptotically more efficient than the naive solution as long as the PRF evaluation is sub-linear in the database size.

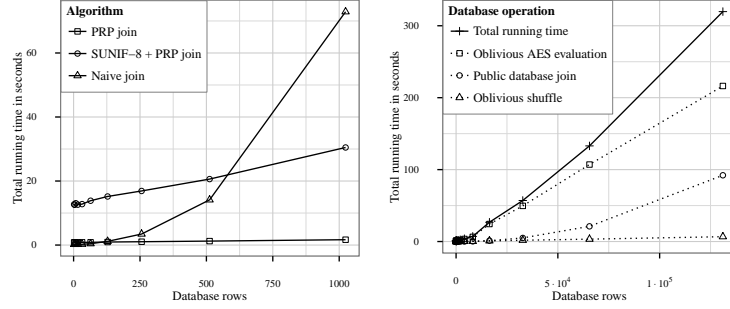


Fig. 2. Benchmarking results for the oblivious database join operation.

hash function and m invocations of OHASH the probability that two different inputs lead to the same output is upper bounded by $\frac{1}{2}m^2\varepsilon$.

Proof (Sketch). The claim about security is evident as multiplication together with addition is sufficient to implement scalar product over \mathbb{F}_2 . The collision probability follows from the union bound $\Pr[\text{collision}] \leq \binom{m}{2} \cdot \varepsilon \leq \frac{m^2\varepsilon}{2}$. \square

Efficiency. A collision in the same key column invalidates the assumptions of Theorem 2, whereas a collision between keys of different tables introduces fraudulent entry in the resulting table. Hence, the size of Carter-Wegman construction must be chosen so that the probability of a collision event is negligible. By using 2^{-80} as the failure probability, we get that 128 bit Carter-Wegman construction allows us to operate up to 33.5 million table entries, which is clearly more than a secure database join protocol can handle in feasible time. To handle around million entries with the same failure probability it is sufficient to use 119-bit Carter-Wegman construction. However, note that the standard implementation of OHASH that computes each bit of the MAC separately and thus duplicates the data vector for each bit, has a larger communication complexity than oblivious AES. Experiments show that for 288 bit input and 128-bit output the complexity of a single OHASH is around 25 ms while the amortised complexity is around 5.7 ms. The corresponding numbers are 11 ms and 0.012 ms for the optimised OHASH protocol detailed in Appendix A. To put the results into context, note that unoptimised OHASH is over 10 times slower than oblivious AES, while the optimised OHASH has almost no impact to performance of the equi-join protocol as its running time is around 5%.

4.3 Benchmarking results

We measured the performance of two secure database join protocols with the same setup as we used for timing the oblivious AES evaluation. For the experiment, we measured how much time it takes to join two database tables consisting

of five 32-bit columns including the single column key. Both databases were of the same size and each key in one table had exactly one matching key in the other table. For AES, we used the BCIRC-1 version of the protocol as it has the lowest amortized cost for tables with thousands of rows.

Results depicted in Figure 2 clearly indicate that PRPJOIN protocols is much more efficient even for modest database sizes and it scales nearly linearly. More precisely, the only non-linear performance component is public database join operation, which is known to take $\Theta(m \log m)$ operations. The exact balance between oblivious AES evaluation and database shuffle depends on the number of columns. As the oblivious database shuffle scales linearly with the number of columns, the fraction of time spent on shuffling increases linearly with the number of columns. However, the slope is rather small.

For instance, consider two database tables with 10,000 rows each. Then the amount of time spent on oblivious shuffle becomes comparable with oblivious AES evaluation only if the number of columns per table exceeds 180 for our experimental setting. Hence, we can safely conclude that the oblivious database join is feasible in practical applications.

The NAIVEJOIN algorithm spends most of its time doing oblivious database shuffle. The shuffle operation itself is efficient, but the share size of the database is big. Even for two tables consisting of 1000 rows we must shuffle a database with million rows. Hence, it is affordable only for small databases.

4.4 Comparison with related work

Protocols for privacy-preserving database join have been proposed before. However, none of them are applicable in our model where input and output tables are secret shared. One of the first articles on privacy-preserving datamining showed how exponentiation can be used to compute equi-join in two-party case [4]. However, their protocol reveals the resulting database.

Freedman *et al.* showed how oblivious polynomial evaluation and balanced hashing can be used to implement secure set intersection [21]. The resulting two-party protocol is based on additively homomorphic encryption and has complexity $\Theta(m_1 m_2)$ without balanced hashing. The latter significantly reduces the amount of computations by splitting the elements into small distinct groups. The same idea is not directly applicable in our setting, since our data is secret shared, while their protocol assumes that key columns are local inputs.

Oblivious polynomial evaluation is not very useful in our context, as it is shorthand for the test $x \in \{b_1, \dots, b_k\}$ which requires $\Theta(k)$ multiplications, while the PRPJOIN protocol does all such comparisons publicly.

Hazay and Lindell [24] have also proposed a similar solution that uses pseudorandom permutation to hide initial data values and performs secure set intersection on ciphertexts. However, they are working in a two-party setting where one of the parties learns the intersection.

5 Conclusion

In this paper we showed that there are several compelling ways to implement oblivious AES evaluation in a multi-party setting where the plaintext and the ciphertext are shared between the parties. As the second important contribution, we described and benchmarked efficient protocols for joining secret-shared databases.

Our benchmarking results showed that it is possible to get throughputs around 3500 blocks per second for the oblivious AES, which is the fastest three-party MPC implementation known to the authors. In general, any block cipher based on substitution permutation networks (SPN) is a good candidate for oblivious evaluation as long as the Sbox has low multiplicative complexity and the rest of the cipher is linear over \mathbb{F}_{2^k} . Experimental results allow us to conclude that throughput around 350 blocks per second is achievable for any comparable SPN cipher, as the evaluation method of [29] is applicable for any Sbox.

Note that the AES key schedule is appropriate for oblivious evaluation, as all the round keys can be computed on demand. Consequently, the usage of pre-shared round keys reduces the running time for a single operation only by 25%. The only way to get more efficient oblivious evaluation protocols is to use Sbox constructions with smaller multiplicative complexity than 32. However, these Sboxes are also more likely to be weaker against linear cryptanalysis and algebraic attacks. Thus, it would be really difficult to come up with more compelling block cipher for multi-party setting – any secure block cipher designed for the oblivious evaluation, is also a good ordinary block cipher.

For the database join, we showed how to combine oblivious evaluation of almost universal hashing and pseudorandom functions to get a collision resistant pseudorandom function, which can handle arbitrary sized database keys. The resulting PRPJOIN protocol works under the assumption that all key column entries are unique. Although we can always fall back to NAIVEJOIN and preserve security without this restriction, the performance penalty is excessive. A better solution remained out of the space restrictions of this paper and is presented in the extended version [1].

From a truly theoretical viewpoint, the question whether sub-quadratic complexity for oblivious database join is achievable depends on existence of pseudorandom functions with low multiplicative complexity. The latter is an interesting open question. Another practically more important open question is to find new almost universal hash function constructions with lower multiplicative complexity or to prove that current constructions are optimal. The circuit complexity of universal hash functions has been studied in the context of energy efficiency [28], the main goal has been minimisation of total circuit complexity which is a considerably different minimisation goal.

References

1. Raknet – multiplayer game network engine. <http://www.jenkinssoftware.com>.
2. SecureSCM. Technical report D9.1: Secure Computation Models and Frameworks. <http://www.securescm.org>, July 2008.

3. Rakesh Agrawal, Alexandre Evfimievski, and Ramakrishnan Srikant. Information sharing across private databases. In *Proceedings of the 2003 ACM SIGMOD'03*, pages 86–97, New York, NY, USA, 2003. ACM.
4. Jean-Philippe Aumasson, Simon Fischer, Shahram Khazaei, Willi Meier, and Christian Rechberger. New Features of Latin Dances: Analysis of Salsa, ChaCha, and Rumba. In *Proceedings of FSE '08*, volume 5086 of *LNCS*, pages 470–488. Springer, 2008.
5. Assaf Ben-David, Noam Nisan, and Benny Pinkas. FairplayMP: a system for secure multi-party computation. In *Proceedings of ACM CCS' 08*, pages 257–266, New York, NY, USA, 2008. ACM.
6. Niv Gilboa Benny Chor and Moni Naor. Private information retrieval by keywords. Cryptology ePrint Archive, Report 1998/003, 1998. <http://eprint.iacr.org/>.
7. D.J. Bernstein. ChaCha, a variant of Salsa20. <http://cr.yp.to/chacha.html>, 2008.
8. Dan Bogdanov, Sven Laur, and Jan Willemsen. Sharemind: A Framework for Fast Privacy-Preserving Computations. In *Proceedings of ESORICS 2008*, volume 5283 of *LNCS*, pages 192–206. Springer, 2008.
9. Dan Bogdanov, Riivo Talviste, and Jan Willemsen. Deploying secure multi-party computation for financial data analysis (Short Paper). In *Proceedings of FC'12*, volume 7397 of *LNCS*, pages 57–64. Springer, 2012.
10. Peter Bogetoft, Dan Christensen, Ivan Damgård, Martin Geisler, Thomas Jakobsen, Mikkel Krøigaard, Janus Nielsen, Jesper Nielsen, Kurt Nielsen, Jakob Pagter, Michael Schwartzbach, and Tomas Toft. Secure multiparty computation goes live. In *Proceedings of Financial Cryptography and Data Security*, volume 5628 of *LNCS*, pages 325–343. Springer, 2009.
11. Joan Boyar and René Peralta. A New Combinational Logic Minimization Technique with Applications to Cryptology. In *Experimental Algorithms*, volume 6049 of *LNCS*, pages 178–189. Springer, 2010.
12. Joan Boyar and René Peralta. A small depth-16 circuit for the aes s-box. In *SEC*, volume 376 of *IFIP Advances in Information and Communication Technology*, pages 287–298. Springer, 2012.
13. Martin Burkhart, Mario Strasser, Dilip Many, and Xenofontas Dimitropoulos. SEPIA: Privacy-preserving aggregation of multi-domain network events and statistics. In *Proceedings of the USENIX Security Symposium '10*, pages 223–239, Washington, DC, USA, 2010.
14. Ran Canetti. Security and composition of multiparty cryptographic protocols. *J. Cryptology*, 13(1):143–202, 2000.
15. Ran Canetti. Universally composable security: A new paradigm for cryptographic protocols. In *proceedings of FOCS '01*, pages 136–145, 2001.
16. Larry Carter and Mark N. Wegman. Universal classes of hash functions. *J. Comput. Syst. Sci.*, 18(2):143–154, 1979.
17. Ivan Damgård, Matthias Fitzi, Eike Kiltz, Jesper Buus Nielsen, and Tomas Toft. Unconditionally secure constant-rounds multi-party computation for equality, comparison, bits and exponentiation. In *Proceedings of TCC*, volume 3876 of *LNCS*, pages 285–304. Springer, 2006.
18. Ivan Damgård and Marcel Keller. Secure multiparty aes. In *Proceedings of Financial Cryptography*, volume 6052 of *LNCS*, pages 367–374. Springer, 2010.
19. Ivan Damgård, Valerio Pastro, Nigel P. Smart, and Sarah Zakarias. Multiparty computation from somewhat homomorphic encryption. In *Proceedings of CRYPTO 2012*, volume 7417 of *LNCS*, pages 643–662. Springer, 2012.

20. Michael Freedman, Kobbi Nissim, and Benny Pinkas. Efficient Private Matching and Set Intersection. In *Proceedings of EUROCRYPT 2004*, volume 3027 of *LNCS*, pages 1–19. Springer, 2004.
21. Craig Gentry and Shai Halevi. Implementing Gentry’s Fully-Homomorphic Encryption Scheme. In *EUROCRYPT 2011*, volume 6632 of *LNCS*, pages 129–148. Springer, 2011.
22. Oded Goldreich. *The Foundations of Cryptography - Volume 2, Basic Applications*. Cambridge University Press, 2004.
23. Carmit Hazay and Yehuda Lindell. Constructions of truly practical secure protocols using standard smartcards. In *ACM Conference on Computer and Communications Security*, pages 491–500, 2008.
24. Wilko Henecka, Stefan Kögl, Ahmad-Reza Sadeghi, Thomas Schneider, and Immo Wehrenberg. TASTY: tool for automating secure two-party computations. In *Proceedings of ACM CCS ’10*, pages 451–462. ACM, 2010.
25. Yan Huang, David Evans, Jonathan Katz, and Lior Malka. Faster Secure Two-Party Computation Using Garbled Circuits. In *Proceedings of 20th USENIX Security Symposium*, pages 8–12, 2011.
26. Yuval Ishai and Anat Paskin. Evaluating branching programs on encrypted data. In *Proceedings of TCC ’07*, volume 4392 of *LNCS*, pages 575–594. Springer, 2007.
27. Jens-Peter Kaps, Kaan Yuksel, and Berk Sunar. Energy scalable universal hashing. *IEEE Trans. Comput.*, 54(12):1484–1495, December 2005.
28. John Launchbury, Iavor S. Diatchki, Thomas DuBuisson, and Andy Adams-Moran. Efficient lookup-table protocol in secure multiparty computation. In *Proceedings of ICFP*, pages 189–200. ACM, 2012.
29. Sven Laur, Riivo Talviste, and Jan Willemson. From oblivious AES to efficient and secure database join in the multiparty setting. Cryptology ePrint Archive, Report 2013/203, 2013. <http://eprint.iacr.org/>.
30. Sven Laur, Jan Willemson, and Bingsheng Zhang. Round-Efficient Oblivious Database Manipulation. In *Proceedings of Information Security*, volume 7001 of *LNCS*, pages 262–277. Springer Berlin / Heidelberg, 2011.
31. Lior Malka. Vmccrypt: modular software architecture for scalable secure computation. In *Proceedings of ACM CCS ’11*, pages 715–724, New York, NY, USA, 2011. ACM.
32. National Institute of Standards and Technology (NIST). Advanced Encryption Standard (AES). *Federal Information Processing Standards Publications*, FIPS-197, 2001.
33. Jesper Buus Nielsen, Peter Sebastian Nordholt, Claudio Orlandi, and Sai Sheshank Burra. A new approach to practical active-secure two-party computation. In *Proceedings of CRYPTO 2012*, volume 7417 of *LNCS*, pages 681–700. Springer, 2012.
34. Benny Pinkas, Thomas Schneider, Nigel Smart, and Stephen Williams. Secure two-party computation is practical. In *Proceedings of ASIACRYPT 2009*, volume 5912 of *LNCS*, pages 250–267. Springer, 2009.

A Carter-Wegman MAC protocol proof

As the computation of Carter-Wegman hash function is essentially a matrix-vector multiplication over the field \mathbb{F}_2 , we can use an optimisation technique,

Input-output specification

Protocol input is a shared s -bit value $[[m]]$ and shared s -bit keys $[[k_1]], \dots, [[k_\ell]]$.
 Protocol output is a shared ℓ -bit MAC value $[[c]]$.

Precomputation phase

1. Each miner \mathcal{P}_i generates ℓ random bits $r_i^1, \dots, r_i^\ell \leftarrow \mathbb{Z}_2$.

Data distribution phase

3. Miner \mathcal{P}_1 sends s -bit shares $m_1, k_{1,1}, \dots, k_{\ell,1}$ to \mathcal{P}_2 .
 Miner \mathcal{P}_2 sends s -bit shares $m_2, k_{1,2}, \dots, k_{\ell,2}$ to \mathcal{P}_3 .
 Miner \mathcal{P}_3 sends s -bit shares $m_3, k_{1,3}, \dots, k_{\ell,3}$ to \mathcal{P}_1 .

Post-processing phase

5. Each miner \mathcal{P}_i computes $w_{ij}^t \leftarrow m_i^t \wedge k_{j,i}^t \oplus m_{i-1}^t \wedge k_{j,i}^t \oplus m_i^t \wedge k_{j,i-1}^t$ for each key $j \in \{1, \dots, \ell\}$ and bit $t \in \{1, \dots, s\}$ and sums them up together with re-randomisation $c_i^j \leftarrow w_{ij}^1 \oplus \dots \oplus w_{ij}^s \oplus r_i^j \oplus r_{i-1}^j$.

Protocol 3: More efficient protocol for Carter-Wegman MAC

which is applicable in many other matrix multiplication settings. The corresponding protocol is depicted as Protocol 3. We use double brackets to denote secret shared values, e.g. the secret shared version of $s = s_1 \oplus s_2 \oplus s_3$ is shown as $[[s]]$, where party \mathcal{P}_i holds s_i . For double indices, the second index shows which party holds the bitstring and the first shows for which output bit it will be used for. Since all values are bitwise shared, we can operate with individual bits of the shares. Operations on individual bits use superscript bit index notation.

Theorem 4. *Assume that the shares of m are correctly generated. Then Protocol 3 is correct and secure against single passively corrupted miner.*

Proof (Sketch). For each bit c^j of MAC the correctness follows from

$$\begin{aligned} [[c^j]] &= \bigoplus_{i=1}^3 \left(\bigoplus_{t=1}^s m_i^t \wedge k_{j,i}^t \oplus m_{i-1}^t \wedge k_{j,i}^t \oplus m_i^t \wedge k_{j,i-1}^t \right) \oplus r_i^j \oplus r_{i-1}^j \\ &= \bigoplus_{i=1}^3 \bigoplus_{t=1}^s (m_i^t \wedge k_{j,i}^t \oplus m_{i-1}^t \wedge k_{j,i}^t \oplus m_i^t \wedge k_{j,i-1}^t) = \bigoplus_{t=1}^s (m^t \wedge k_j^t) = h(k_j, m) \end{aligned}$$

since the inner most sum contains all combinations of $m_a \wedge k_b$.

For the security analysis, it is sufficient to consider the corruption of \mathcal{P}_2 who receives all shares owned by \mathcal{P}_1 . Note that two shares out of three have always uniform distribution. Hence, it is trivial to simulate all messages received by \mathcal{P}_2 . Since \mathcal{P}_2 is semihonest, the simulator can extract shares of the message and keys from the input of \mathcal{P}_2 and submit them to the trusted party who will return shares c_2^1, \dots, c_2^ℓ . Since the simulator knows what random values r_2^1, \dots, r_2^ℓ \mathcal{P}_2 is going to use, it can pick r_1^1, \dots, r_1^ℓ so that \mathcal{P}_2 will indeed output c_2^1, \dots, c_2^ℓ . We leave the detailed analysis of the simulation construction to the reader. \square