

Differentially Oblivious Database Joins: Overcoming the Worst-Case Curse of Fully Oblivious Algorithms

Shumo Chu* Danyang Zhuo Elaine Shi T-H. Hubert Chan
UCSB Duke CMU HKU

shumo@cs.ucsb.edu, danyang@cs.duke.edu, runting@cs.cmu.edu, hubert@cs.hku.hk

Abstract

Numerous high-profile works have shown that access patterns to even encrypted databases can leak secret information and sometimes even lead to reconstruction of the entire database. To thwart access pattern leakage, the literature has focused on *oblivious* algorithms, where obliviousness requires that the access patterns leak nothing about the input data.

In this paper, we consider the `Join` operator, an important database primitive that has been extensively studied and optimized. Unfortunately, any *fully oblivious Join* algorithm would require *always* padding the result to the *worst-case* length which is *quadratic* in the data size N . In comparison, an insecure baseline incurs only $O(R + N)$ cost where R is the true result length, and in the common case in practice, R is relatively short. As a typical example, when $R = O(N)$, any fully oblivious algorithm must inherently incur a prohibitive, N -fold slowdown relative to the insecure baseline. Indeed, the (non-private) database and algorithms literature invariably focuses on studying the *instance-specific* rather than *worst-case* performance of database algorithms. Unfortunately, the stringent notion of full obliviousness precludes the design of efficient algorithms with non-trivial instance-specific performance.

To overcome this worst-case performance barrier of full obliviousness and enable algorithms with good instance-specific performance, we consider a relaxed notion of access pattern privacy called (ϵ, δ) -differential obliviousness (DO), originally proposed in the seminal work of Chan et al. (SODA'19). Rather than insisting that the access patterns leak no information whatsoever, the relaxed DO notion requires that the access patterns satisfy (ϵ, δ) -differential privacy. We show that by adopting the relaxed DO notion, we can obtain efficient database `Join` mechanisms whose instance-specific performance *approximately matches* the insecure baseline, while still offering a meaningful notion of privacy to individual users. Complementing our upper bound results, we also prove new lower bounds regarding the performance of any DO `Join` algorithm.

Differential obliviousness (DO) is a new notion and is a relatively unexplored territory. Following the pioneering investigations by Chan et al. and others, our work is among the very first to formally explore how DO can help overcome the worst-case performance curse of full obliviousness; moreover, we motivate our work with database applications. Our work shows new evidence why DO might be a promising notion, and opens up several exciting future directions.

*Author order is randomized.

1 Introduction

We consider a scenario in which a trusted client (e.g., an Intel SGX enclave or a trusted client laptop) outsources an encrypted database to an untrusted storage provider (e.g., untrusted memory or a cloud server). The client would like to make queries to the database without endangering the privacy of users in the database. Since all data contents are encrypted, the key challenge is to ensure that *access patterns* to the database do not accidentally harm individual users’ privacy. Notably, plenty of recent attacks [IKK12, CGPR15, NKW15, GLMP19, KKNO16, GMN⁺16, KPT20] against encrypted database systems (e.g., CryptDB [PRZB11], Cipherbase [ABE⁺13], and TrustedDB [BS14]) showed that when left unprotected, access patterns can leak highly sensitive information, and in some cases, even lead to reconstruction of the entire database. Given the importance of this problem, several high-profile works implemented *oblivious* database systems, including Opaque [ZDB⁺17], OblIDB [EZ19], Obladi [CBC⁺18], as well as the work by Arasu and Kaushik [AK14], where *obliviousness* requires that access patterns leak nothing about the underlying data.

In this paper, we focus on an important database operation, the **Join** operation, which has been studied extensively in the database literature [Yan81, Sha86, AHV95, GUW09, NNRR14, CBS15, KNRR15, BKS17, HY19, AGM08]. Given two tables and a specified attribute (henceforth also called the *join key* or *key* for short)¹, the **Join** operation computes, for each possible join key value k , the Cartesian product of the rows in each table with the join key k . For example, if the join key k appears twice in the first table and three times in the second table, then in the join result the join key k will have six occurrences. Some works focus on the special case of foreign-key join where it is promised that in one of the input tables, each key appears only once. In this paper, however, we consider the more general case where a key may have multiple occurrences in both tables.

Unfortunately, existing oblivious database systems [ZDB⁺17, EZ19, AK14] do not provide a satisfactory solution for the **Join** operation. A fundamental problem is that any *fully oblivious* algorithm for **Join** must always incur the *worst-case cost* even when the actual join result may be short. To see this, recall that an algorithm is said to be *oblivious* iff its memory access patterns (and runtime²) are indistinguishable for any two inputs of the same length [GO96, Gol87, SCSL11]. For the case of **Join**, the result size for the worst-case input is $\Theta(N_1 \cdot N_2)$ where N_1 and N_2 denote the sizes of the two input tables, respectively. Thus any fully oblivious algorithm must incur at least $\Omega(N_1 \cdot N_2)$ cost on *any* input instance, even when the input instance has a short join result³. This is very expensive in real-world databases, since the join result is typically much smaller than quadratic in the common case.

In this paper, we ask the following natural question:

Can we design join algorithms that provide a meaningful and mathematically rigorous notion of privacy, and moreover, avoid having to pay the worst-case quadratic penalty on every input?

Parametrized algorithm design and instance-specific performance. We follow the well-established paradigm in the algorithms and database literature, and focus on *instance-specific* complexity measures. *Parametrized analysis and instance-specific performance have been widely adopted in the database, algorithms, as well as cryptography literature*, and its importance has been explained in numerous prior works. For example, a line of work in the classical (non-private) database literature focused on optimizing the instance-specific performance for database

¹A join key can also be a set of attributes, without loss of generality, we assume two tables are joined on single attribute in this paper.

²Note that the length of the physical accesses is the same as the program’s runtime.

³The naïve solution of simulating an insecure join algorithm with Oblivious RAM does not provide full obliviousness unless the runtime is padded to the worst case.

joins [HY19, BPWZ14, Yan81, KNRR15, NNRR14] — in this context the output length is often used as an additional parameter to characterize the algorithms’ performance. The study of instance-specific performance is also closely related to the prominent line of work on parametrized algorithms design and analysis [Rou20, MY15, CKX06, CFK⁺15, FG06, SO92] (see Roughgarden’s textbook for an excellent overview [Rou20].) Last but not the least, notable works in the cryptography literature also consider how to achieve good instance-specific performance (sometimes called “input-specific runtime” in the cryptography literature) in the Turing Machine or the Random Access Machine (RAM) models: for example, the seminal work by Goldwasser et al. [GKP⁺13] is motivated by “overcoming the worst-case curse” in cryptographic constructions; and a similar notion is adopted in subsequent works [KP16, AFS19].

Prior approaches introduce arbitrary leakages to achieve good instance-specific performance. As mentioned, full obliviousness *precludes* the design of algorithms with non-trivial instance-specific performance. However, prior oblivious database systems [ZDB⁺17, EZ19, AK14] do care about instance-specific performance bounds. To achieve good instance-specific performance, they give up on full obliviousness (even though this line of work is commonly referred to as “oblivious databases”), and introduce arbitrary leakages, e.g., by leaking the multiplicity of keys, the lengths of intermediate arrays, the final output length, and/or the exact runtime of the (non-private) program. The ramifications of such leakages are poorly understood, and can lead to unforeseen privacy breaches. Since numerous prior encrypted database systems allowing arbitrary leakages have been broken [IKK12, CGPR15, NKW15, GLMP19, KKNO16, GMN⁺16, KPT20], our philosophy is to advocate for an approach that provides rigorous mathematical guarantees on the leakage.

Although our work focuses on a privately outsourced database scenario, it is interesting to note that in the cryptography literature, a line of work has focused on general RAM computations on encrypted data [GKK⁺12, GHL⁺14, GHRW14, CCC⁺16, BCP15a, CCHR16]. These works also care about plugging access pattern leakage. Thus, to achieve full security, essentially full obliviousness is necessary in these constructions (and indeed these constructions rely on Oblivious RAM as a building block). Typically this line of work either pads the RAM computation to the runtime on the worst-case input, or they allow leakage of the exact runtime and thus violate full obliviousness — the ramification of such leakage is unclear and can lead to severe privacy breaches in some applications.

Overcoming the worst-case curse with differential obliviousness. Since the worst-case performance curse is inherent for full obliviousness, we would need a relaxed (but nonetheless meaningful and rigorous) privacy notion to achieve good instance-specific performance. We therefore turn our attention to the notion of *differential obliviousness* (DO) recently defined by Chan et al. [CCMS19]. Simply put, differential obliviousness requires that the access patterns revealed during a program’s execution must satisfy (ϵ, δ) -differential privacy [DMNS06]. So far, a couple of prior works [CCMS19, BNZ19] have shown theoretical separations between DO and full obliviousness, thus providing initial theoretical evidence why DO is worth studying. Besides the few pioneering investigations, the landscape of DO remains much unexplored. Designing DO algorithms, especially for practically motivated applications, is a relatively new territory.

1.1 Our Contributions and Results

We show novel DO database join algorithms that can *approximately match the instance-specific performance of the insecure baseline*, whereas any fully oblivious join algorithm must inherently incur, on common instances with $O(N)$ -sized outputs, at least a linear (in database size) blowup

relative to the insecure baseline. Our work is among the very first to formally explore how DO can help overcome the worst-case curse of full obliviousness.

Specifically, we present *two main upper bound results*: 1) a DO join algorithm in the standard word-RAM model where the primary performance metrics are the algorithm’s runtime and output length, and 2) a DO join algorithm in the *external-memory* model where the primary performance metrics are the algorithm’s cache complexity and output length. Both algorithms approximately match the performance of the insecure baselines in the corresponding setting. Both models are important to consider: the standard word-RAM model is the prevalent model in which algorithms are studied; and the external-memory model is the best fit when we rely on secure processors such as Intel’s SGX to privately outsource the sensitive database to an untrusted server (as we explain more later).

We also prove *lower bound* results regarding the performance of any DO join algorithm. The lower bounds show that some small slowdown relative to the insecure baseline is necessary. Moreover, our upper bound matches the lower bound when the result size is at least quasi-linear. For other parameter regimes, e.g., when the result size is linear or shorter, there remains a small gap between our upper bounds and lower bounds — and bridging this gap is an interesting direction for future work.

We now present the result statements more formally.

Results for the word-RAM model. Recall the application scenario mentioned at the beginning of the paper: a *trusted* client stores an *encrypted* database on an *untrusted* storage. Data can only be decrypted within the trusted client which also runs the database engine. Anything fetched or written to the storage is encrypted such that the adversary can only observe the access patterns.

In the standard word-RAM model, we assume that the trusted client is a CPU with $O(1)$ private registers, and the cost is measured in terms of the number of memory words transmitted between the CPU and memory (which equates to the runtime of the algorithm). Table 1 summarizes our results for the standard RAM model. For simplicity, the results are stated for the typical parameters⁴ $\epsilon = \Theta(1)$ and $\delta = 1/N^c$ for some constant $c \geq 1$ and a more generalize version will be provided in Theorem 1.1.

As shown in Table 1, our algorithm achieves $O(R + N \log N)$ runtime and $R + O((\mu_{\max} + \log N) \cdot \log N)$ result size where N denotes the total input length, R denotes the true result size (when the insecure algorithm is run), and μ_{\max} denotes the multiplicity of the most frequent join key in the database. Note that even an insecure join algorithm must incur at least $R + N$ runtime since it has to at least read the input and write down the output. In the common case in practice, the true output size R is small, e.g., $R = O(N)$. In this case, our DO algorithm achieves almost a factor of N performance improvement relative to any fully oblivious solution whose cost is inherently quadratic.

We also compare our algorithm with a naïve DO algorithm that basically simulates the insecure algorithm (described in Section 3.5) using the state-of-the-art *statistically* secure Oblivious RAM [CS17, WCS15]⁵ and then appends an appropriate noise to the result as well as to the program’s runtime (see Section 3.5 for a more detailed description). In comparison, our algorithm is a $\log^2 N$ factor faster than the naïve DO algorithm. Other standard DO techniques such as the work by Komargodski and Shi [KS21] fail to work in our context — see Section 2.5 for more discussions. In light of our lower bound for any DO join algorithm, our upper bound achieves optimality in terms

⁴In the cryptography literature, sometimes we want δ to be a negligible function in N . In this case, the $\log N$ factor in the performance bound is replaced with any super-logarithmic function.

⁵Like Chan et al. [CCMS19], we adopt a statistical notion of differential obliviousness that defends against even computationally unbounded adversaries.

Table 1: Our results: stated for the typical parameters when $\epsilon = \Theta(1)$ and $\delta = \frac{1}{N^c}$ for some constant $c \geq 1$. N_1 and N_2 denote the lengths of the two input tables, $N := N_1 + N_2$, R denotes the length of the true join result, and μ_{\max} denotes the maximum multiplicity of any join key in the two input tables. $\Theta(\cdot)$ means that it is both an upper- and lower-bound.

	Runtime	Result size
Insecure	$\Theta(R + N)$	R
Fully oblivious	$\Theta(N_1 \cdot N_2)$	$\Theta(N_1 \cdot N_2)$
Our differentially oblivious algorithms		
Naïve: Theorem 3.11	$O((R + N \log N) \log^2 N)$	$R + O(N \log N)$
Main scheme 1: Thm 1.1	$O(R + N \log N)$	$R + O((\mu_{\max} + \log N) \cdot \log N)$
LB: Thm 1.2	$\Omega(R + N \log \log N + \mu_{\max} \log N)$	$R + \Omega(\mu_{\max} \log N)$

of result length as long as $\mu_{\max} \geq \log N$, and is optimal in terms of runtime when $\mu_{\max} = \Theta(N)$.

Our main theorems⁶ are also informally described below for a broader range of choices for ϵ and δ than Table 1.

Theorem 1.1 (Our DO join algorithm). *Let R be the length of the true join result (i.e., without fillers), let μ_{\max} denote the multiplicity of the most frequent join key in either input array, and let N denote the total input length. There is an (ϵ, δ) -differentially oblivious join algorithm that runs in time $O(R + N(\log N + \frac{1}{\epsilon} \log \frac{1}{\delta}))$ and produces a result whose length is at most $R + O(\frac{1}{\epsilon} \cdot (\mu_{\max} + \frac{1}{\epsilon} \log \frac{1}{\delta}) \cdot \log \frac{1}{\delta})$.*

Table 1 also shows our lower bound which states that any DO join algorithm must incur at least $\Omega(R + N \log \log N + \mu_{\max} \log N)$ runtime and must have a result size of at least $R + \Omega(\mu_{\max} \log N)$ with high probability (assuming the same typical choices of ϵ and δ). A formal statement with more general parameters is given below.

Theorem 1.2 (Limits of any DO join algorithm (informal)). *Let N be the total input length, then, for most reasonable choices⁷ of ϵ and δ ,*

1. *any (ϵ, δ) -differentially oblivious join algorithm must produce a result of at least $R + \Omega(\mu_{\max} \cdot \frac{1}{\epsilon} \cdot \log \frac{\epsilon}{\delta})$ with at least δ/ϵ probability.*
2. *any “natural” (ϵ, δ) -differentially oblivious join algorithm must have some input of total length N and whose true join result size is R , such that the algorithm incurs at least $\Omega(R + N \log \log \frac{1}{\delta} + \mu_{\max} \cdot \frac{1}{\epsilon} \cdot \log \frac{\epsilon}{\delta})$ runtime with at least δ/ϵ probability.*

In the above, the lower bound for runtime holds for a broad class of *natural* algorithms that do not perform encoding or computation on the elements’ payloads — indeed, most known join algorithms fall into this class. We refer the reader to Section 6.2 for more details.

⁶In Table 1, we assume that $\delta = 1/N^c$ like the standard differentially privacy literature suggests. In some cryptographic application settings where one may desire δ to be a negligible function in N , there will be an extra (arbitrarily small) super-constant factor added to the bounds in Table 1. See also Theorem 1.1 for the statement with general parameters.

⁷See Theorem 6.4 for a more precise characterization of the parameter regime in which the lower bound holds.

Table 2: Our results: cache-agnostic cache complexity. See the caption of Table 1 for the meaning of the notations N_1 , N_2 , N , and R . Our results need to assume the standard “tall cache” and “wide block” assumptions, i.e., $M \geq B^2$, and $B \geq \log^{0.55} N$ where M is the cache size and B is the block size.

Cache-oblivious cache complexity	
Insecure	$O(\frac{R}{B} + \frac{N}{B} \cdot \log_{\frac{M}{B}} \frac{N}{B})$
Fully oblivious	$\Theta(N_1 \cdot N_2/B)$
Our differentially oblivious algorithms	
Naïve: Theorem 3.11	$O((R + N \log N) \log N \cdot \log_B N)$
Main scheme 2: Cor 1.3	$O(\frac{R}{B} + \frac{N}{B} \cdot \log N)$

Results for in the cache-agnostic, external-memory model. The external-memory model [AV88, Vit01, FLPR99] and cache complexity are important for scenarios where we want to employ secure processors to enable encrypted, differentially oblivious databases. Imagine that the server has a secure processor such as Intel SGX. In this case, the database is stored in an encrypted format on the server, and only the secure processor can decrypt the data and perform computation. In other words, we can think of the *trusted client* as the secure processor, and the rest of the server’s software stack is untrusted. Moreover, a remote client can communicate with the secure processor using a secure channel to ask queries and receive answers back.

Interestingly, it turns out that when Intel SGX is used to outsource both the computation and storage to an untrusted server, the major performance metric is the number of pages the SGX enclave needs to fetch. Each enclave page swap is a heavy-weight operation that involves communication with the untrusted operating system, and moreover, the enclave must decrypt (or encrypt) the memory page being swapped in (or out). In this scenario, the trusted enclave memory can be viewed as a *cache* whose size is henceforth denoted M , and each page is a *block* (i.e., the atomic unit being swapped in and out) whose size is henceforth denoted B . Further, the rest of the storage outside the trusted enclave memory is the *external memory*. An algorithm’s *cache complexity* is defined as the number of blocks transmitted between the cache and the external memory during the algorithm’s execution. In Appendix B.1, we provide additional background on the external-memory model which is a well-accepted model in the algorithms literature — it is very interesting to observe that the line of work on external-memory algorithms [AV88, Vit01, FLPR99] is a perfect fit for studying the performance of algorithms running on commodity secure processors.

We propose a variant of our algorithm optimized for cache complexity. Our algorithm is *cache agnostic*, i.e., the algorithm is unaware of the cache’s parameters, namely, M and B . Cache-agnostic was also commonly referred to as “cache-oblivious” in the algorithms literature [FLPR99, Dem02]. In our paper, we use the term “cache-agnostic” instead to disambiguate from our usage of the term “obliviousness”. The importance and advantages of cache-agnostic algorithms have been extensively discussed in the algorithms literature [FLPR99, Dem02]. First, a cache-agnostic algorithm is “universal” and the performance bounds hold no matter what the system parameters (including M and B) are. Not only so, when deployed on a multi-level memory hierarchy, an optimal cache-agnostic algorithm would give optimal IO performance between any two adjacent levels of the hierarchy [FLPR99, Dem02].

Table 2 summarizes our cache complexity results. The insecure baseline and the naïve DO algorithm in the table are described in Section 3.5. As shown in the table, our cache complexity

is quite close to that of the insecure baseline, and outperforms the naïve DO algorithm by a $(B/\log B) \cdot \log^2 N$ factor. In a typical scenario, $B = \text{poly log } N$; in this case, our improvement over the naïve DO algorithm is polylogarithmic.

Last but not the least, our cache efficient instantiation has relatively small constants in the big-O notation, and therefore an interesting future direction is to implement our algorithm and measure its concrete efficiency. We summarize our cache-complexity results in the following corollary:

Corollary 1.3 (Our DO join algorithm: cache complexity). *There is an (ϵ, δ) -DO database join algorithm that incurs cache complexity upper bounded by $\frac{1}{B} \cdot O\left(N\left(\log \frac{M}{B} \frac{N}{B} + \frac{1}{\epsilon} \log \frac{1}{\delta}\right) + R + \left(\frac{1}{\epsilon} \log \frac{1}{\delta}\right)^2\right)$, assuming the standard tall cache assumption $M = \Omega(B^2)$ and the wide block assumption $B = \Omega(\log^{0.55} N)$.*

Both the tall cache assumption and the wide block assumption are standard assumptions adopted commonly in the external-memory algorithms line of work [AV88, Vit01, FLPR99, Dem02, ABD⁺07].

Technical highlight. Inspired by the original work of Chan et al. [CCMS19], we adopt the following design paradigm for devising DO algorithms. At a very high level, we decompose the task of designing a DO algorithm into the following: 1) identify a set of intermediate ideal functionalities with *differentially private leakage*; and 2) leverage oblivious algorithm building blocks to obliviously realize these ideal functionalities, such that the access patterns leak only the stated differentially private leakage, and nothing else.

Although the design paradigm is simple to state, the non-trivial challenge is to identify appropriate intermediate functionalities that not only lend to solving our problem, but also being cognizant that the computational tasks they embody must have efficient oblivious realizations. We defer the algorithmic details to subsequent formal sections, and we hope that our algorithmic techniques can inspire the design of DO algorithms for new applications. We believe that our work provides further evidence on top of the early-stage explorations of Chan et al. [CCMS19] and Beimel et al. [BNZ19] that differential obliviousness is a useful notion that deserves attention.

2 Technical Roadmap

For convenience, henceforth we call the two input tables *arrays*, denoted \mathbf{I}_1 and \mathbf{I}_2 respectively. Each element in \mathbf{I}_1 and \mathbf{I}_2 is either a real element of the form (k, v) or a filler element of the form (\perp, \perp) . For a real element, k is called the join key (or key for short) and v is called the payload. The database join operation wants to compute, for each unique join key k , the Cartesian product of the elements contained in both arrays with join key k . All results are concatenated and output, and moreover, the output is allowed to contain an arbitrary number of filler elements that may be needed for privacy.

Remark 2.1 (Simplifying assumption for the roadmap). Throughout our informal technical roadmap, we will assume the typical parameters $\epsilon = \Theta(1)$ and $\delta = 1/N^c$ for an arbitrary constant $c \geq 1$. In this case, $\frac{1}{\epsilon} \log \frac{1}{\delta} = \Theta(\log N)$, and we thus use two expressions interchangeably — but our formal sections later will differentiate the two to be more general. Specifically, jumping ahead, we often need to add noises of magnitude roughly $\frac{1}{\epsilon} \log \frac{1}{\delta} = \Theta(\log N)$.

Our differential oblivious join algorithm will make use of standard oblivious algorithm building blocks including oblivious sorting [AKS83, ACN⁺20, RS21], and oblivious compaction [AKL⁺20a, LSX19]. We review these building blocks in more detail in Section 3.6.

2.1 Warmup Algorithm

We first present a warmup that achieves $O(R + N \log^2 N)$ runtime — the warmup algorithm does not achieve the bounds stated earlier; but it is conceptually simpler and helps our understanding. Later in Section 2.3, we describe additional techniques to improve the algorithm’s asymptotical performance, and achieve the bounds stated earlier.

A strawman idea. To understand our algorithm, let us first consider a flawed strawman — we sketch the high-level idea, and for the time being, omit the details on how to oblivious sorts to implement the relevant steps.

1. *Compute the bin load array L .* First, using a constant number of oblivious sorts, write down a list L of length $N := |\mathbf{I}_1| + |\mathbf{I}_2|$. Each element in L is of the form $(k, \hat{n}_k^{(1)}, \hat{n}_k^{(2)})$ where $\hat{n}_k^{(b)}$ denotes the noisy count of the join key k in table $b \in \{1, 2\}$. The noisy count is obtained by adding an appropriate, independently sampled noise to the actual multiplicity of join key k in the corresponding array. To maintain correctness, the noise must be non-negative. So rather than adding a Laplacian noise with standard deviation $1/\epsilon$, we shift the Laplacian to the right to be centered at $U/2 = \Theta(\frac{1}{\epsilon} \cdot \log \frac{1}{\delta}) = \Theta(\log N)$. In this way, the noise lies within the range $[0, U]$ except with δ probability⁸.

The list L should contain all join keys that appear in at least one input array, padded with fillers of the form $(\star_1, \hat{n}_{\star_1}^{(1)}, \hat{n}_{\star_1}^{(2)})$, $(\star_2, \hat{n}_{\star_2}^{(1)}, \hat{n}_{\star_2}^{(2)})$, \dots , to a length of N . The filler join keys \star_1, \star_2, \dots have an actual multiplicity of 0 in both input arrays, and thus their noisy counts are the shifted Laplacian noise in the range $[0, U]$. The list L is sorted by the join key k , and all filler join keys appear at the end.

2. *Binning.* Now, we have $2N$ bins each indexed by a pair (b, i) where $b \in \{1, 2\}$ and $i \in [N]$. Let k_i denote the i -th smallest join key. Then, the bin indexed (b, i) has capacity $\hat{n}_{k_i}^{(b)}$, and all elements in \mathbf{I}_b with the join key k_i are destined for this bin. Using a constant number of oblivious sorts, route all elements in either array to their respective destined bins, and pad each bin with fillers to its intended capacity. Note that the bins corresponding to the filler join keys \star_1, \star_2, \dots have no real elements in them and are full of fillers.
3. *Bin-wise Cartesian product.* Now, take every pair of bins $(1, i)$ and $(2, i)$ for $i \in [N]$, and compute the Cartesian product of elements in the two bins — if the two elements being joined have the same real join key, add the joined tuple to the output; otherwise, add a filler element to the output.

A flaw that violates differential obliviousness. The above algorithm is natural and conceptually simple; it almost works, except for a critical flaw that violates differential obliviousness, which is illustrated in Figure 1 and explained in detail below. Observe that the array L is sorted by the join key during Step 1. Now, consider an input $\mathcal{I} := (\mathbf{I}_1, \mathbf{I}_2)$ in which 8th smallest join key k_8 appears only 1 time in \mathbf{I}_1 and does not appear in \mathbf{I}_2 , and all other join keys that appear in \mathcal{I} appear more than $16U$ times in both arrays. Consider a 2-neighboring input where the only occurrence of k_8 is replaced with a join key $k_{-\infty}$ that 1) does not exist in \mathcal{I} , and 2) is smaller than all other join keys in \mathcal{I} . In this case, an adversary observing the access patterns of the program can easily tell which input is used. Recall that the bin pairs are sorted in the order of the join keys: in the case

⁸In our formal sections later, we actually use a shifted geometric distribution which is the discrete counterpart of the real-valued Laplacian, and moreover we simply truncate the δ -probability mass outside the range $[0, U]$ which allows us to get deterministic bounds on the algorithm’s runtime.

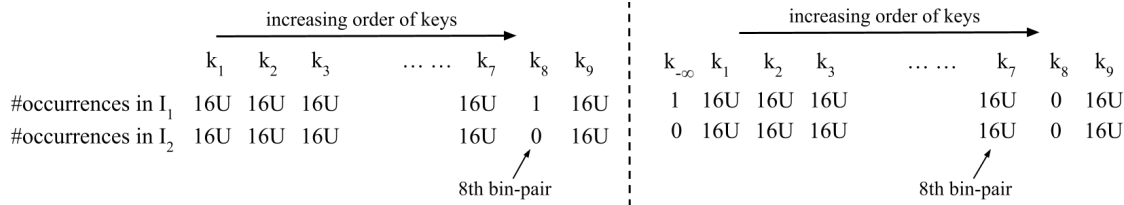


Figure 1: Strawman scheme: a flaw that violates differential obliviousness. See the main body for more explanation.

of \mathcal{I} , the 8-th bin pair has small capacities and the first bin pair has large capacities (where small means between at most $U + 1$ and large means at least $16U$). In the case of \mathcal{I}' , however, the first bin pair, now corresponding to the join key $k_{-\infty}$, has small capacities.

A remedy. It turns out that a simple fix can address the above flaw: instead of ordering the array L using the join key k , we can order it lexicographically based on the $(\hat{n}_k^{(1)}, \hat{n}_k^{(2)})$ fields. Intuitively, this can avoid accidental information leakage through the ordering. More specifically, in the above steps 2 and 3, the capacities of all bins are leaked to the adversary. Therefore, if we use the fields $(\hat{n}_k^{(1)}, \hat{n}_k^{(2)})$ to order the array L and the corresponding bins, then the only information leaked is the *multiset* of $(\hat{n}_k^{(1)}, \hat{n}_k^{(2)})$ pairs. Given the *multiset* of the $(\hat{n}_k^{(1)}, \hat{n}_k^{(2)})$ pairs, the access patterns of the above algorithm are fully determined.

Therefore, it suffices to prove that the leakage, i.e., the *multiset* of the $(\hat{n}_k^{(1)}, \hat{n}_k^{(2)})$ pairs, satisfies $(O(\epsilon), O(\delta))$ -differential privacy. In our formal technical sections later, we shall prove that this is indeed the case as long as the noises are chosen from the shifted and truncated geometric distribution defined in Section 3.4.

Final step: compaction of the join result. The above modification fixes the security flaw, but at this moment, the result output by our algorithm may have length $O(R + N \log^2 N)$ — see Section 2.2 for a more detailed analysis. Specifically, when the true result length R is small, the additive term $N \log^2 N$ is dominant.

We would like our algorithm to output a result that is as short as possible specific to the instance. Clearly, for the result to be correct, it cannot be shorter than R . In Section 6, we prove that all (ϵ, δ) -DO algorithms must have result length at least $R + \Omega(\mu_{\max} \cdot \frac{1}{\epsilon} \log \frac{1}{\delta})$ where μ_{\max} is the maximum multiplicity of any join key in either array.

Our idea is to obviously compact the result output by the above algorithm to length $R + \text{noise}$ where noise is sampled from an appropriate distribution. The most natural idea is to sample a shifted Laplacian noise proportional to $\Delta_{\text{global}}/\epsilon$ where Δ_{global} is the global sensitivity of the exact result length (i.e., how much the length of the exact result would change in the *worst case* when we change one position in the input). However, Δ_{global} can be as large as $\Theta(N)$. Instead, we would like to achieve an *instance-optimal* bound on the result length (for almost all parameter regimes). To do so, we add noise proportional to the local sensitivity (or instance-specific sensitivity) which is equal to μ_{\max} . To make the idea work, however, we have to first obtain a noisy version $\hat{\mu}_{\max}$ to the local sensitivity μ_{\max} , and then add noise proportional to $\hat{\mu}_{\max}$ to the result length. In this way, our algorithm achieves result length $R + (\mu_{\max} + \log N) \log N$ which is instance-optimal in light of the $R + \Omega(\mu_{\max} \log N)$ lower bound in almost all parameter regimes. We defer a detailed description of the scheme to the formal technical sections.

2.2 Performance of the Warmup Algorithm

As mentioned, the warmup algorithm, obtained by changing the way the array L is ordered in the strawman solution, achieves runtime $O(R + N \log^2 N)$. To understand the techniques in Section 2.3 that improves the performance to $O(R + N \log N)$, let us first understand the performance breakdown.

1. The first step, which computes the bin load array L , performs a constant number of oblivious sorts on arrays of length at most $N := |\mathbf{I}_1| + |\mathbf{I}_2|$, and thus takes $N \log N$ time.
2. The second step, which places the elements into bins, performs a constant number of oblivious sorts, and the length of the arrays sorted is upper bounded by the sum of the bin capacities. In the worst case, there can be $\Theta(N)$ sparsely loaded bins each with only $O(1)$ number of real elements. The noise added to each bin's capacity is roughly of magnitude $\frac{1}{\epsilon} \log \frac{1}{\delta}$ which is $O(\log N)$ under typical parameters (see Remark 2.1). Therefore, the length of the array sorted is at most $O(N \log N)$, and the second step takes time $O((N \log N) \log(N \log N)) = N \log^2 N$.
3. The third step computes the Cartesian product of pairs of bins: the runtime of this step is the actual result length R when there is no noise, plus the number of fillers. The number of fillers is maximized when there are $\Theta(N)$ bins each with $O(1)$ number of real elements and $O(\log N)$ fillers. In this case, the total number of fillers after the Cartesian product is $O(N \log^2 N)$. Therefore, the third step takes time $O(R + N \log^2 N)$.

Summarizing the above, our warmup algorithm achieves $O(R + N \log^2 N)$ runtime.

2.3 Final Algorithm

Section 2.2 reveals that the $N \log^2 N$ additive term in the performance bound is incurred because in the worst-case scenario, there can be $\Theta(N)$ sparsely loaded bin-pairs each with $O(1)$ real elements, and padded with $\Theta(\log N)$ fillers⁹. This introduces the $N \log^2 N$ additive term in two ways: 1) the binning step requires sorting arrays of length $O(N \log N)$ which takes $O(N \log^2 N)$ time; and 2) the bin-wise Cartesian product introduces $O(N \log^2 N)$ fillers.

Imprecisely speaking, having many sparsely loaded bin-pairs cause a small “signal to noise” ratio, i.e., the ratio of fillers is high. To improve the performance bound to $O(R + N \log N)$, our idea is to reduce the number of bin-pairs to $O(N/\log N)$, thereby improving the “signal to noise” ratio. We say that a join key k is *sparse* iff its noise counts $\widehat{n}_k^{(1)}$ and $\widehat{n}_k^{(2)}$ are both upper bounded by $2U$, where recall that $U = \Theta(\frac{1}{\epsilon} \log \frac{1}{\delta})$. A join key that is not sparse is said to be *dense*. Recall that $N := |\mathbf{I}_1| + |\mathbf{I}_2|$.

As mentioned, our noise distribution is upper bounded by U except with probability δ . This means that if a bin-pair contains a dense join key, then at least one of the bins in the pair has at least U real elements in it except with δ probability. Thus there can be no more than $O(N/\log N)$ bin-pairs for dense-keys. Our focus therefore is to consolidate multiple sparse join keys into the same bin-pair such that each bin-pair contains at least $\Theta(U)$ elements (including elements from both input arrays). To achieve this, we perform the following.

Compute key-to-bin mapping. Recall that earlier we computed the bin load array L which contains tuples of the form $(k, \widehat{n}_k^{(1)}, \widehat{n}_k^{(2)})$ sorted according to lexicographical ordering on $(\widehat{n}_k^{(1)}, \widehat{n}_k^{(2)})$.

⁹For example, this can happen if there are many elements with $O(1)$ occurrences in both input arrays, or with $O(1)$ occurrences in one array but not appearing in the other. In the latter case, essentially there are many elements in the symmetric difference of the two input arrays that do not contribute to the true joined result.

It is not too hard to extend the algorithm for computing L such that the bin load array L also stores the actual counts, i.e., L now contains entries of the form $(k, n_k^{(1)}, \hat{n}_k^{(1)}, n_k^{(2)}, \hat{n}_k^{(2)})$. where $n_k^{(1)}$ and $n_k^{(2)}$ denote the actual multiplicity of the join key k in \mathbf{I}_1 and \mathbf{I}_2 , respectively.

Now, we can classify L into a part L_s corresponding to sparse keys, i.e, $L_s := \{(k, n_k^{(1)}, \hat{n}_k^{(1)}, n_k^{(2)}, \hat{n}_k^{(2)}) \in L : \hat{n}_k^{(1)} \leq 2U, \hat{n}_k^{(2)} \leq 2U\}$; and a part $L_d := L \setminus L_s$ corresponding to dense join keys. All of L, L_s , and L_d are sorted according to lexicographical ordering on $(\hat{n}_k^{(1)}, \hat{n}_k^{(2)})$; and L_s and L_d can be constructed in $O(N)$ time if we allow the access patterns to reveal the noisy counts contained in L .

Our goal now is to construct an array called **BinMap** that maps join keys to bin pairs (note by constructing **BinMap**, we have not moved the elements into their bins yet — the actual moving will be done in the subsequent binning step). Each entry of **BinMap** is of the form (k, i) , meaning that the join key k should be mapped to the i -th bin-pair. To construct **BinMap**, we first scan through L_d : for each $i \in \{1, 2, \dots, |L_d|\}$, if the i -th entry in L_d has the join key k , add the tuple (k, i) to the array **BinMap**. At this moment, **BinMap** stores the mapping from each dense join key to its destined bin-pair index. The capacities of these bin-pairs (for dense join keys) are determined by the noisy counts in L_d .

Next, we will add to **BinMap** the mapping from sparse join keys to their bins. Specifically, sparse join keys are mapped to additional bin-pairs numbered $\{(1, j), (2, j) : j = |L_d| + 1, |L_d| + 2, \dots, |L_d| + O(N/U)\}$ where j is also called the bin-pair index. All of these bins (for sparse join keys) will have capacity *exactly* $4U$, and here we allow multiple join keys to be mapped to the same bin pair. The invariant we want is that for each bin-pair $(1, j), (2, j)$ where $j \in \{|L_d| + 1, |L_d| + 2, \dots, |L_d| + O(N/U)\}$, a total of at least $2U$ elements will be mapped to the bin pair (summing across both input arrays). In this way, all the sparse join keys altogether will not consume more than $N/2U$ bins.

To achieve this, we can scan linearly through L_s . For each entry $(k, n_k^{(1)}, \hat{n}_k^{(1)}, n_k^{(2)}, \hat{n}_k^{(2)}) \in L_s$ encountered during the scan, we append to **BinMap** a tuple (k, j) that indicates that join key k is mapped to the j -th bin-pair. Here j is the current bin-pair counter whose starting value is $|L_d| + 1$, and j is incremented whenever one of the current bins is about to exceed its capacity. To achieve this, the algorithm additionally maintains two counters that remember how many cumulative elements have been mapped to the current bins $(1, j)$ and $(2, j)$ so far. Whenever one of the bins $(1, j)$ or $(2, j)$ is about to exceed its capacity $4U$, we increment the bin-pair counter j and reset both counters to be 0 again, i.e., we start mapping join keys to the next bin-pair. We stress that the access pattern of this step is fixed and depends only on $|L_s|$ because we append a single entry to **BinMap** whenever we visit an entry of L_s .

Binning. Next, we perform a binning step and move elements into their desired bins — henceforth a bin designated for a single dense join key is called a D-bin, and a bin designated for possibly multiple sparse join keys is called an S-bin. To perform the binning, we need **BinMap** which provides the mapping between join keys and their bin indices. D-bins have capacities determined by the corresponding noisy counts contained in L_d and there are $|L_d|$ D-bins. S-bins have capacities exactly $4U$, and the number of S-bins is an a-priori fixed upper bound CN/U for a sufficiently large constant C . We can now use a constant number of oblivious sorts to move elements into their desired bins, padding each bin with fillers to its intended capacity as defined above. Note that it is possible that some S-bins do not receive any element.

Remainder of the algorithm. The remainder would be similar to the warmup algorithm. We perform bin-wise Cartesian product; and finally we obviously compact the result adding an appropriate noise to the final result length. Since now, multiple join keys can share the same bin-pair,

during the Cartesian product step, whenever we try to pairwise-join two elements with different join keys, we append a filler to the output array.

2.4 Lower Bound Results

As mentioned, we prove new lower bounds on the result length and runtime of any DO join algorithm. The result length lower bound is proven using the definition of differential obliviousness. Our lower bound on runtime is obtained by taking the maximum of two lower bounds: 1) the aforementioned lower bound on the result length, and 2) a lower bound that stems from a privacy-preserving and complexity-preserving reduction from sorting to database join. Such a reduction shows that any DO join algorithm must suffer from the same lower bound for DO sorting proven recently by Chan et al. [CCMS19]. We refer the reader to Section 6. for the detailed statements and proofs.

2.5 Additional Related Work

In this paper, we adopt the differential obliviousness notion defined by Chan et al. [CCMS19]. Besides Chan et al. [CCMS19], several other works also considered related but somewhat incomparable notions [MG18, KKNO17, WCM18].

Besides the aforementioned works on *oblivious* databases [ZDB⁺17, EZ19, AK14, CBC⁺18], a related but incomparable line of work [PRZB11, CJJ⁺13, JJK⁺13, PRV⁺11, CGKO06, SWP00, SBC⁺07, AHKM19, SPS14, CJJ⁺14, BS14] focuses on encrypted databases or searchable encryption systems. Typically, this line of works either give up on hiding access patterns [PRZB11, CJJ⁺14, SPS14, CJJ⁺13, CGKO06, SWP00, SBC⁺07, BS14], or hide the access pattern by performing a linear scan through the entire database upon every update or query [AHKM19]. The cryptographic techniques for computation on encrypted data developed in this line of work is somewhat orthogonal and complementary to our techniques for obfuscating the access patterns.

Following the classical differential privacy (DP) literature, another line of work that focuses on differentially private database queries [KTM⁺19, CSS11, HR10, GRU12, KTH⁺19]. These works assume that the database curator is trusted and only aims to guarantee that the result is DP — they are not concerned about information leakage through the runtime behavior of the database engine. Notably, the techniques we use to release the noisy counts for the multiplicity of join keys may be remotely reminiscent of differentially private histogram mechanisms [BV18, KKMN09, BNS16, AHKM19, Sur19, BDB16]. We stress, however, that our work and techniques are of a different nature from classical DP mechanisms, including classical DP algorithms for releasing histograms. While prior DP mechanisms introduce noise to the statistics released, in our context, we introduce noise to the algorithm’s access patterns (and not its output) — importantly, we need to do so without affecting the algorithm’s correctness. It would be very interesting, however, to apply our DO techniques to classical DP algorithms — in this way, both the runtime behavior of the database as well as the released statistics would guarantee DP, i.e., we get “end-to-end” privacy.

Mazloom and Gordon [MG18] proposed new techniques that guarantee differential obliviousness for tasks that can be performed in a graph-parallel framework. They rely on shuffling to guarantee that the access patterns of these algorithms reveal only differentially private histograms. While seemingly related to our techniques, we do not know any straightforward way to apply their techniques to the join problem and get our asymptotical bounds: partly, our techniques are non-trivial because *we avoid suffering too much overhead for join keys that are in the symmetric difference of the two input arrays* — these join keys do not contribute to the true joined result. Imprecisely speaking, doing such “pruning” *privately* introduces non-trivial algorithmic challenges.

Komargodski and Shi [KS21] suggest how to compile any Turing Machine (TM) to a differentially oblivious TM. A strawman idea is to apply their compiler to the (insecure) TM that computes the database join problem. Unfortunately, this completely fails because their work defines neighboring on the *operational sequences* of two TMs; whereas we define neighboring on the *inputs*. For two inputs that are neighboring (i.e., Hamming distance 1), applying the insecure TM that computes database join over these inputs may result in operational sequences that are far apart. This is also partly why our problem is challenging.

2.6 Open Problems

Our work is among the first to explore DO algorithms motivated by practical database systems. Our work reveals that this is a promising direction with many intriguing open questions. For example, can we bridge the gap between our upper- and lower-bounds for a broader range of parameters?

Another exciting direction is to explore DO algorithms for other common database queries. In this paper, we considered two-way joins. In the classical, non-private database literature, however, *multi-way join* [Yan81, Sha86, AHV95, GUW09, NNRR14, CBS15, KNRR15, BKS17, HY19, AGM08] received significantly more attention because instance-optimal two-way join is long known to be a solved problem. Our paper reveals that with the extra privacy requirements, even two-way join raises non-trivial algorithmic challenges. Of course, it also makes sense to ask whether one can design efficient DO algorithms for multi-way joins as well, especially, whether we can (approximately) match the performance of the best known insecure algorithms. Recent works in the database literature also considered join algorithms for the special case when the input tables are already sorted based on the join key (or more generally, preprocessed in some way). In this case, it may not be necessary to read the entire input, and sublinear (non-private) algorithms are known [NNRR14, KNRR15]. Therefore, an open question is whether we can achieve sublinear DO algorithms for sorted inputs. Besides joins, more general class of database queries such as conjunctive queries [AHV95] are also interesting to consider.

Last but not the least, as mentioned, the cache-efficient variants of our algorithms are potentially implementable and suitable for SGX-type scenarios. Implementing DO algorithms in practical database systems and evaluating their concrete performance is another exciting future direction.

3 Preliminaries

We assume that the algorithm is executed in a standard Random Access Machine (RAM) model. The adversary can observe the access patterns of the program, i.e., in each step, which memory location is accessed and whether each access is a read or write operation. The adversary, however, cannot observe the data contents — for example, in a secure processor setting, the data contents protected by encryption. We will be concerned about two metrics: 1) the program’s runtime; and 2) the program’s cache complexity [AV88]. For the runtime statements, we assume a standard word-RAM and the CPU has only $O(1)$ number of private registers. For the cache complexity metric, we assume a standard external-memory RAM model [AV88] where the CPU has M bits of private cache, and every time it needs to transmit data to and from memory, an atomic unit (called a “block”) of B bits is transferred. The cache complexity metric measures how many blocks are transmitted between the CPU and memory. All of our algorithms are cache agnostic [FLPR99, Dem02], i.e., the algorithm need not know the cache’s parameters M and B .

3.1 Database Join

Database join is the following problem. Let \mathcal{K} denote the space of join key and \mathcal{V} denote the space of payloads. Let \mathbf{I}_1 and \mathbf{I}_2 be two input arrays each containing pairs of the form (k, v) , where $k \in \mathcal{K} \cup \{\perp\}$ is called a *join key* and $v \in \mathcal{V} \cup \{\perp\}$ is called the *payload*. If $k \in \mathcal{K}$ and $v \in \mathcal{V}$, the element (k, v) is said to be a *real* element. Without loss of generality, we may assume that if an element's join key k is \perp , its payload v must be \perp too: such elements, of the form (\perp, \perp) , are said to be *filler* elements.

Our goal is to output an array \mathbf{O} such that for each non-filler join key $k \in \mathcal{K}$ that appears in both \mathbf{I}_1 and \mathbf{I}_2 : let $\{(k, v_1), (k, v_2), \dots, (k, v_m)\}$ be the multi-set of elements having join key k in \mathbf{I}_1 , and let $\{(k, w_1), (k, w_2), \dots, (k, w_{m'})\}$ be the multi-set of elements having join key k in \mathbf{I}_2 ; then the multi-set $\{k, v_i, w_j\}_{i \in [m], j \in [m']} \subseteq \mathbf{O}$. We use R to denote the size of this multi-set. Moreover, besides the multi-set $\{k, v_i, w_j\}_{i \in [m], j \in [m']}$, \mathbf{O} should contain no other element with the join key k . Additionally, the output array \mathbf{O} may contain any number of filler elements of the form (\perp, \perp, \perp) .

In other words, the output array \mathbf{O} contains, for each join key $k \in \mathcal{K}$, the Cartesian product of the elements in both input arrays under join key k ; and additionally, \mathbf{O} may contain some filler elements. The filler elements in the output array \mathbf{O} may be needed for privacy reasons as will become clear later.

Remark 3.1 (Motivation for allowing fillers in the input array). In our formulation, we allow the input arrays \mathbf{I}_1 and \mathbf{I}_2 to contain filler elements, because the length of the input arrays may already be noisy to mask the true number of elements contained in it. For example, if the input array comes from a differentially oblivious database such as in the work by Chan et al. [CCMS19], then the input arrays would already contain a random number of filler elements.

3.2 Full Obliviousness

In this paper, we consider execution of algorithms on the Random Access Machine (RAM) model. Let Alg denote a possibly randomized algorithm and let \mathbf{I} denote an input to the algorithm. We use the notation $\text{Accesses}^{\text{Alg}}(\mathbf{I})$, a random variable denoting the sequence of memory addresses accessed and whether each access is a read or write, generated by a random execution of the algorithm Alg on input \mathbf{I} . Therefore, $\text{Accesses}^{\text{Alg}}(\mathbf{I})$ is also called the “access patterns” of Alg on input \mathbf{I} .

Definition 3.2 (δ -obliviousness). *We say that an algorithm Alg satisfies δ -obliviousness w.r.t. the leakage function $\text{Leak}(\cdot)$, iff there exists a simulator Sim , such that $\text{Accesses}^{\text{Alg}}(\mathbf{I})$ has statistical distance at most δ from the simulated access patterns $\text{Sim}(\text{Leak}(\mathbf{I}))$.*

In other words, the access patterns are simulatable by a simulator Sim which knows only the leakage function but nothing more about the input \mathbf{I} . Note also that δ is allowed to be a function in $N = |\mathbf{I}|$.

A typical leakage function is leaking only the length of the input and nothing else, i.e., $\text{Leak}(\mathbf{I}) := |\mathbf{I}|$.

Definition 3.3 (Full obliviousness). *Henceforth, whenever we say Alg is (fully) oblivious (i.e., omitting the leakage function and δ), the leakage function would be the default one $\text{Leak}(\mathbf{I}) := |\mathbf{I}|$, and δ is assumed to be a negligible function in N .*

Throughout the paper, we say that a function $\nu(N)$ is a negligible function, iff for any $c \in \mathbb{N}$, there exists a sufficiently large N_0 such that for all $N \geq N_0$, $\nu(N) \leq 1/N^c$. In other words, ν drops faster than any inverse-polynomial function.

3.3 Differential Obliviousness

Neighboring inputs. Two inputs $(\mathbf{I}_1, \mathbf{I}_2)$ and $(\mathbf{J}_1, \mathbf{J}_2)$ are said to be neighboring, iff $|\mathbf{I}_1| = |\mathbf{J}_1|$ and $|\mathbf{I}_2| = |\mathbf{J}_2|$, and moreover, the following holds:

- either $\mathbf{I}_1 = \mathbf{J}_1$, and moreover, \mathbf{I}_2 and \mathbf{J}_2 differ in exactly one position;
- or $\mathbf{I}_2 = \mathbf{J}_2$, and moreover, \mathbf{I}_1 and \mathbf{J}_1 differ in exactly one position.

Differential obliviousness. Imagine that a database join algorithm Alg is executed on a Random Access Machine (RAM). The two input arrays $(\mathbf{I}_1, \mathbf{I}_2)$ reside in memory, and at the end of the algorithm, the output array \mathbf{O} is written to a designated position in memory.

Definition 3.4 ((ϵ, δ) -differential obliviousness). *We say that a database join algorithm Alg satisfies (ϵ, δ) -differential obliviousness or (ϵ, δ) -DO for short, iff for any neighboring inputs $(\mathbf{I}_1, \mathbf{I}_2)$ and $(\mathbf{J}_1, \mathbf{J}_2)$, for any set S ,*

$$\Pr \left[\text{Accesses}^{\text{Alg}}(\mathbf{I}_1, \mathbf{I}_2) \in S \right] \leq e^\epsilon \cdot \Pr \left[\text{Accesses}^{\text{Alg}}(\mathbf{J}_1, \mathbf{J}_2) \in S \right] + \delta$$

where $\text{Accesses}^{\text{Alg}}(\mathbf{I}_1, \mathbf{I}_2)$ is a random variable denoting the sequence of memory addresses (also called access patterns) generated by a random execution of the algorithm Alg on input $(\mathbf{I}_1, \mathbf{I}_2)$.

Specifically, in the standard RAM model, in each time step, the machine visits one memory location, reading it and then updating it with either the old value or a new value. Therefore, $\text{Accesses}^{\text{Alg}}(\mathbf{I}_1, \mathbf{I}_2)$ is just the ordered list of all memory addresses accessed in all time steps. Moreover, the length of $\text{Accesses}^{\text{Alg}}(\mathbf{I}_1, \mathbf{I}_2)$ is also the (randomized) runtime of the algorithm. Like the standard notion of differential privacy, our notion secures against unbounded adversaries.

Typical choices of ϵ and δ . Typically, we would like $\epsilon = \Theta(1)$. The standard differential privacy literature [Vad17] recommends that δ be set to $1/N^c$ for some constant $c > 1$. In the cryptography literature, sometimes we would like δ to be negligibly small in N .

3.4 Mathematical Building Blocks

Definition 3.5 (Symmetric geometric distribution). *Let $\alpha > 1$. The symmetric geometric distribution $\text{Geom}(\alpha)$ takes integer values such that the probability mass function at k is $\frac{\alpha-1}{\alpha+1} \cdot \alpha^{-|k|}$.*

As we shall see, our algorithm will hide the true cardinality of a set by padding it with a random number of filler elements. Below we define a useful distribution from which we shall sample the noises.

Definition 3.6 (Shifted and truncated geometric distribution). *Let $\epsilon > 0$ and $\delta \in (0, 1)$ and $\Delta \geq 1$. Let k_0 be the smallest positive integer such that $\Pr[|\text{Geom}(e^{\frac{\epsilon}{\Delta}})| \geq k_0] \leq \delta$, where $k_0 = \frac{\Delta}{\epsilon} \ln \frac{2}{\delta} + O(1)$. The shifted and truncated geometric distribution $\mathcal{G}(\epsilon, \delta, \Delta)$ has support in $[0, 2(k_0 + \Delta - 1)]$, and is defined as:*

$$\min\{\max\{0, k_0 + \Delta - 1 + \text{Geom}(e^\epsilon)\}, 2(k_0 + \Delta - 1)\}$$

For the special case $\Delta = 1$, we write $\mathcal{G}(\epsilon, \delta) := \mathcal{G}(\epsilon, \delta, 1)$.

In the main body of the paper, for simplicity, we shall first assume that we can sample from the shifted and truncated geometric distribution in $O(1)$ time. Later in Appendix B, we discuss how to remove this assumption without blowing up the runtime.

Notation. Given two random variables X and Y , we use $X \sim_{(\epsilon, \delta)} Y$ to denote that X and Y satisfy the standard (ϵ, δ) -differentially private inequality, i.e., for all subsets S , $\Pr[X \in S] \leq e^\epsilon \cdot \Pr[Y \in S] + \delta$, and $\Pr[Y \in S] \leq e^\epsilon \cdot \Pr[X \in S] + \delta$,

Fact 3.7 (Differential privacy through adding truncated and shifted geometric noise [CCMS19, BV18]). *Let $\epsilon > 0$ and $\delta \in (0, 1)$. Suppose u and v are two non-negative integers such that $|u - v| \leq \Delta$. Then,*

$$u + \mathcal{G}(\epsilon, \delta, \Delta) \sim_{(\epsilon, \delta)} v + \mathcal{G}(\epsilon, \delta, \Delta).$$

Fact 3.8 (Post-processing). *Let $X \in \mathcal{X}$ and $X' \in \mathcal{X}$ be random variables and let $F : \mathcal{X} \rightarrow \mathcal{Y}$ be a possibly randomized function. Suppose that $X \sim_{(\epsilon, \delta)} X'$. Then, we have that*

$$F(X) \sim_{(\epsilon, \delta)} F(X').$$

Fact 3.9 (Composition of differentially private mechanisms (Theorem B.1 of [DR14])). *Suppose that for any neighboring \mathbf{I} and \mathbf{I}' , $T_1(\mathbf{I}) \sim_{(\epsilon_1, \delta_1)} T_1(\mathbf{I}')$. Henceforth let $\text{supp}(T_1(\mathbf{I}))$ denote the support of applying the function T_1 to the data \mathbf{I} . Suppose that for any neighboring \mathbf{I} and \mathbf{I}' , for any $s_1 \in \text{supp}(T_1(\mathbf{I})) \cup \text{supp}(T_1(\mathbf{I}'))$, $T_2(\mathbf{I}, s_1) \sim_{(\epsilon_2, \delta_2)} T_2(\mathbf{I}', s_1)$. Then, for any neighboring \mathbf{I}, \mathbf{I}' , we have that*

$$(T_2, T_1)(\mathbf{I}) \sim_{(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)} (T_2, T_1)(\mathbf{I}')$$

The following operational lemma will guide our algorithm design and analysis.

Lemma 3.10 (Operational lemma for differential obliviousness [CCMS19]). *If an algorithm is δ_0 -oblivious w.r.t. a leakage function $\text{Leak}(\cdot)$, and moreover, Leak is (ϵ, δ) -differentially private, then the algorithm satisfies $(\epsilon, \delta_0 + \delta)$ -differential obliviousness.*

3.5 Naïve Solutions

Insecure algorithm. If privacy is not needed, there is an algorithm that computes the join result in (expected) $O(N + R)$ time, where R is the actual result size and $N := |\mathbf{I}_1| + |\mathbf{I}_2|$. Basically, hash elements in each input array into a separate Cuckoo hash table [PR04, ANS09] and within each entry of the hash table, use a linked list to store all elements with the same join key. Now, we can pairwise-join elements with the same join key from the two input arrays by querying the Cuckoo hash table. The entire algorithm completes in $O(N + R)$ expected time (where the randomness comes from hashing).

The Cuckoo hashing based insecure solution, however, does not have great cache complexity. For cache complexity, we use a different insecure baseline, that is, first sort the two input arrays by join key, and then use the most straightforward approach to compute pairwise-join elements with the same join key from the two arrays. This algorithm achieves $O(\frac{R}{B} + \frac{N}{B} \cdot \log \frac{M}{B})$.

Observe also that any insecure algorithm must at least read the entire input to guarantee correctness. Thus an always correct insecure algorithm must incur at least $\Omega(N)$ runtime.

Fully oblivious algorithm. The following naïve algorithm, which tries every potential pair of elements from the two input arrays, can achieve fully oblivious database join with $O(|\mathbf{I}_1| \cdot |\mathbf{I}_2|)$ runtime:

For each element $(k, v) \in \mathbf{I}_1$, for each element $(k', v') \in \mathbf{I}_2$: if $k = k'$, add (k, v, v') to the output; else add (\perp, \perp, \perp) to the output.

As argued in Section 6, $O(|\mathbf{I}_1| \cdot |\mathbf{I}_2|)$ is also the best we can hope for if full obliviousness is desired.

Naïve differentially oblivious algorithm. A naïve approach to achieve differentially oblivious join is to rely on a *statistically* secure Oblivious RAM (ORAM) that defends against unbounded adversaries, such as Circuit ORAM [WCS15, CS17], to simulate the aforementioned insecure algorithm with $O(\log^2 N)$ blowup in runtime. Suppose that the result size is R and let N be the total input length. We know that the insecure algorithm must complete in $T \leq C \cdot (R + N)$ steps for some constant C , and it produces an output of length R . Now, given the $O(\log^2 N)$ ORAM simulation overhead, simulating the insecure algorithm with ORAM requires at most $C'(R + N) \log^2 N$ steps for some sufficiently large $C' > C$. However, if we just stopped here, then the algorithm would leak information through the result length R and its running time (which is the same as the total length of the physical memory accesses in the RAM model).

To plug this leakage, the algorithm proceeds to add noise to the result length and its own runtime. To achieve this, we continue to use the ORAM to simulate the following steps:

1. Append $\xi := \mathcal{G}(\epsilon, \delta, N) = O(\frac{N}{\epsilon} \log \frac{1}{\delta})$ number of filler elements to the joined result — this requires the ORAM to simulate $O(\xi)$ additional steps. Henceforth let $\widehat{R} := R + \xi$.
2. Finally, pad the running time of the simulated algorithm to $C'(\widehat{R} + N) \log^2 N$.

Note that in the above, we add noise $\mathcal{G}(\epsilon, \delta, N)$ noise to R because the global sensitivity of R is upper bounded by N , that is, changing one element in the input can change R by at most N .

To obtain better cache complexity, we can place the ORAM schemes' binary tree data structures in an Emde Boas layout; specifically, each access in the original program will incur a cache complexity of $O(\log N \cdot \log_B N)$ in the ORAM simulation.

Theorem 3.11 (Naïve DO algorithm). *The above naïve algorithm satisfies $(\epsilon, \delta + \text{negl}(N))$ -differential obliviousness where $\text{negl}(\cdot)$ denotes a suitable negligible function.*

Further, suppose that $\epsilon = \Theta(1)$ and $\delta = \frac{1}{\text{poly}(N)}$, let $U = O(\frac{1}{\epsilon} \log \frac{1}{\delta}) = O(\log N)$; then, the above naïve algorithm achieves $O((R + NU) \log^2 N)$ runtime and $O((R + NU) \log N \cdot \log_B N)$ cache complexity, and outputs a result of $O(R + NU)$ length.

Proof. The runtime, cache complexity and output length follows from the above description. Observe that because of ORAM (which can fail with $\text{negl}(N)$ probability), the access pattern are simulatable given N and \widehat{R} . By Lemma 3.10, it suffices to prove that the leakage \widehat{R} satisfies (ϵ, δ) -differential privacy.

Consider changing one element in the input $(\mathbf{I}_1, \mathbf{I}_2)$ to form $(\mathbf{I}'_1, \mathbf{I}'_2)$. We have that $|R(\mathbf{I}_1, \mathbf{I}_2) - R(\mathbf{I}'_1, \mathbf{I}'_2)| \leq N$ where $R(\mathbf{I}_1, \mathbf{I}_2)$ denotes the length of the exact result on the input $(\mathbf{I}_1, \mathbf{I}_2)$. By Fact 3.7, we have that

$$\widehat{R}(\mathbf{I}_1, \mathbf{I}_2) \sim_{(\epsilon, \delta)} \widehat{R}(\mathbf{I}'_1, \mathbf{I}'_2).$$

where $\widehat{R}(\mathbf{I}_1, \mathbf{I}_2)$ denotes the length of the result output by the naïve DO algorithm upon input $(\mathbf{I}_1, \mathbf{I}_2)$. Hence, the algorithm is $(\epsilon, \delta + \text{negl}(N))$ -differentially oblivious, where the extra $\text{negl}(N)$ comes from the failure probability of ORAM. \square

3.6 Oblivious Algorithm Building Blocks

We describe several oblivious algorithm building blocks. Unless otherwise noted, obliviousness is defined w.r.t. the input-length leakage (i.e., informally, speaking, only the input length is leaked).

Oblivious compaction. Given an input array where some elements are marked as distinguished, output an array where all distinguished elements are moved to the front, and all non-distinguished

elements are moved to the end. The very recent works by Asharov et al. [AKL⁺20a, AKL⁺20b] constructed a $O(n)$ -time oblivious compaction algorithm that can compact any input array of length n . Their linear-time compaction algorithm is not *stable*, i.e., among the distinguished (or non-distinguished) elements, the output does not preserve the relative order the elements appeared in the input, and this non-stability is inherent [LSX19].

To get our cache complexity result, we will adopt the randomized, cache-agnostic, oblivious compaction algorithm by Lin, Shi, and Xie [LSX19]: their algorithm achieves optimal $O(n/B)$ cache complexity and $O(n \log \log n)$ runtime assuming that $M = \Omega(B^2)$ and $B \geq \log^{0.55} n$.

Oblivious sort. Ajtai, Komlós, and Szemerédi [AKS83] showed that there is a sorting circuit with $O(n \log n)$ comparators that can correctly sort any input array containing n elements. Such a sorting circuit can be executed on a Random Access Machine (RAM) in $O(n \log n)$ time assuming that each element can be represented using $O(1)$ words.

The recent work by Ramachandran and Shi [RS20] constructed a *randomized*, cache-agnostic, oblivious sort algorithm that achieves $O(n \log n)$ runtime and $O((n/B) \log_{M/B}(n/B))$ cache complexity, assuming the tall cache assumption that $M = \Omega(B^2)$ and further $M = \Omega(\log^{1+\epsilon} n)$ for an arbitrarily small constant $\epsilon \in (0, 1)$.

Given oblivious sorting, we can realize a couple intermediate abstractions including oblivious send-receive and oblivious bin placement which we define below. Both primitives can be realized by invoking oblivious sorting constant number of times.

Oblivious send-receive. The send-receive primitive¹⁰ solves the following problem. In the input, there is a source array and a destination array. The source array represents n senders, each of whom holds a key and a value; it is promised that all join keys are distinct. The destination array represents n' receivers each holding a join key. Now, have each receiver learn the value corresponding to the join key it is requesting from one of the sources. If the join key is not found, the receiver should receive \perp . Note that although each receiver wants only one value, a sender can send its values to multiple receivers.

Prior works [CS17, CCS17, BCP15b] have shown that oblivious send-receive can be accomplished through a constant number of oblivious sorts on arrays of length $O(n + n')$. The algorithm is oblivious w.r.t. the leakage n and n' , (i.e., informally, only the lengths of the input arrays are leaked).

Oblivious bin placement. A bin placement algorithm solves the following problem. Suppose we have m bins each of capacity s_1, s_2, \dots, s_m , respectively. We are given an input array denoted \mathbf{I} , where each element is either a *filler* denoted \perp or a *real* element that is tagged with a bin identifier $\beta \in [m]$ denoting which bin it wants to go to. It is promised that every bin will receive no more elements than its capacity. Now, move each real element in \mathbf{I} to its desired bin. If any bin is not full after the real elements have been placed, pad it with filler elements at the end to its desired capacity. Finally, output the concatenation of the resulting bins.

Chan and Shi [CS17] describes an oblivious bin placement algorithm that solves the special case of the problem when all the bin sizes are equal. Their algorithm relies on a constant number of oblivious sorts. It is not difficult to extend their algorithm to the case when the bin sizes are not equal. For completeness, we describe the modified algorithm in Appendix A.1. Specifically, *the algorithm satisfies obliviousness w.r.t. to the leakage that contains the input size, as well as the sizes of all bins.* Let n denote the size of the input array, and let $S := \sum_{\beta \in [m]} s_\beta$ be the sum

¹⁰The send-receive abstraction is often referred to as oblivious routing in the data-oblivious algorithms literature [BCP15b, CS17, CCS17]. We avoid the name “routing” because of its other connotations in the algorithms literature.

of the sizes of all bins. The runtime of the algorithm is upper bounded by $T_{\text{sort}}(n + S)$ and the cache-agnostic, cache complexity of the algorithm is upper bounded by $Q_{\text{sort}}(n + S)$, where $T_{\text{sort}}(n')$ and $Q_{\text{sort}}(n')$ denote the runtime and (cache-agnostic) cache complexity of oblivious sort over an input array of size n' .

4 Warmup Algorithm

We describe our warmup algorithm in detail and give a formal analysis. The intuitions of the warmup algorithm were explained earlier in Section 2.1. Therefore we jump directly into the detailed construction.

4.1 Step 1: Compute Bin Loads

Henceforth let $\{\star_1, \star_2, \dots\}$ be a special, imaginary join keys that do not exist in either input array. In the first step of the algorithm, we would like to compute a bin load array L of length exactly $|\mathbf{I}_1| + |\mathbf{I}_2|$, containing the following:

- Let $\text{Keys}(\mathbf{I}_1 \cup \mathbf{I}_2) \subseteq \mathcal{K}$ denote the union of non-filler join keys in the union of the two input arrays. For each $k \in \text{Keys}(\mathbf{I}_1 \cup \mathbf{I}_2)$, there is an entry of the form $(k, \hat{n}_k^{(1)}, \hat{n}_k^{(2)})$ in L where $\hat{n}_k^{(1)} \in \mathbb{N} \cup \{0\}$ denotes the noisy multiplicity of the join key k in \mathbf{I}_1 , and $\hat{n}_k^{(2)} \in \mathbb{N} \cup \{0\}$ denotes the noisy multiplicity of the join key k in \mathbf{I}_2 .
- Besides the above, L should be padded with entries of the form $(\star_1, \hat{n}_{\star_1}^{(1)}, \hat{n}_{\star_1}^{(2)})$, $(\star_2, \hat{n}_{\star_2}^{(1)}, \hat{n}_{\star_2}^{(2)})$, $(\star_3, \hat{n}_{\star_3}^{(1)}, \hat{n}_{\star_3}^{(2)})$, \dots , to a capacity of $|\mathbf{I}_1| + |\mathbf{I}_2|$. Since the special join keys $\{\star_1, \star_2, \dots\}$ do not exist in either input array, their noisy counts are just noisy approximations of the true multiplicity 0.

Jumping ahead, later on in our algorithm, we will place each input array \mathbf{I}_1 and \mathbf{I}_2 into bins based on their join keys, and the bins' capacities will be determined by the counts, the $\hat{n}_k^{(1)}$'s and $\hat{n}_k^{(2)}$'s. These counts will therefore be leaked to the adversary, but the actual join keys will remain hidden from the adversary. We employ the following algorithm to compute L .

Step 1: Computing the bin load array L

- Obliviously sort $\mathbf{I}_1 \parallel \mathbf{I}_2$ by the elements' join keys. For the same join key, break ties arbitrarily. In the sorted array, the last element for a particular join key k is said to be the representative for k . In a linear scan with two additional registers to keep track of the count of the current join key: 1) each representative element for a join key $k \in \mathcal{K}$ writes down the true multiplicity of k in both \mathbf{I}_1 and \mathbf{I}_2 ; whereas each non-representative element tags itself as a filler.
- Obliviously sort the resulting array, moving all fillers to the end. The front of the array now contains entries of the form $(k, n_k^{(1)}, n_k^{(2)})$ where $k \in \text{Keys}(\mathbf{I}_1 \cup \mathbf{I}_2)$, and $n_k^{(1)}$ and $n_k^{(2)}$ denote the true multiplicity of the join key k in \mathbf{I}_1 and \mathbf{I}_2 , respectively.
- In a linear scan, replace all the filler entries at the end with special entries of the form $(\star_1, 0, 0)$, $(\star_2, 0, 0)$, $(\star_3, 0, 0)$, and so on.
- In a linear scan, add to every count a freshly sampled noise^a from $\mathcal{G}(\epsilon, \delta)$.

- (e) Finally, sort obviously this array based on lexicographical ordering of the fields $(\hat{n}_k^{(1)}, \hat{n}_k^{(2)})$, and output the resulting array as L .

^aLater in the final algorithm in Section 5, we also need to save the original true counts here.

The access patterns of Step 1 leak nothing beyond the lengths of the input arrays, $|\mathbf{I}_1|$ and $|\mathbf{I}_2|$. Further the oblivious sorting dominates the runtime of Step 1. We formally state these observations in the following fact.

Fact 4.1. *Step 1 is oblivious w.r.t. to the leakage $|\mathbf{I}_1|$ and $|\mathbf{I}_2|$. Moreover, the running time is $O(N \log N)$, where $N = |\mathbf{I}_1| + |\mathbf{I}_2|$.*

Proof. The obliviousness claim follows directly from the observation that Step 1 only performs sequential scan and oblivious sorting on arrays of length $N = |\mathbf{I}_1| + |\mathbf{I}_2|$. The runtime is dominated by the oblivious sorting and thus upper bounded by $O(N \log N)$. \square

4.2 Step 2: Binning

In the second step, we would like to put elements in \mathbf{I}_1 and \mathbf{I}_2 into bins based on their join keys. The bin load array L computed earlier determines which bin each element is placed into, and each bin is then padded with filler elements to the intended capacity specified by L . For example, suppose that $L[i] = (k, \hat{n}_k^{(1)}, \hat{n}_k^{(2)})$. Then for $b \in \{1, 2\}$, all elements in \mathbf{I}_b with the join key k will be placed in the bin indexed (b, i) , and this bin will be padded to a capacity of $\hat{n}_k^{(b)}$.

To accomplish this bin placement, we can simply employ the following algorithm. Henceforth we may assume that each element in L is tagged with its index within the array L , since this can be accomplished with a linear scan.

Step 2: Binning

- (a) Invoke an oblivious send/receive algorithm (see Section 3.6) such that each real element in \mathbf{I}_1 and \mathbf{I}_2 receives from L which bin it is destined for.
- (b) Invoke an oblivious bin placement algorithm (see Section 3.6) to place all real elements in \mathbf{I}_1 and \mathbf{I}_2 to its destined bin; and moreover, all bins are padded with fillers to its capacity prescribed by the L array.

Note that the access pattern of Step 2 reveals the boundary between the bins, i.e., the capacities of all bins are revealed; further, it leaks the lengths of \mathbf{I}_1 and \mathbf{I}_2 . These exactly correspond to the leakage of the oblivious bin placement algorithm we employ. However, beyond the bin capacities and the lengths of \mathbf{I}_1 and \mathbf{I}_2 , nothing additional is leaked by the access patterns of Step 2. Henceforth we use the notation RevealM to denote the capacities of all bins. RevealM can be obtained from the array L but removing the join key field from each entry. We therefore have the following fact.

Fact 4.2. *Let RevealM be the same as the array L but removing the join key field from each entry. Step 2 is oblivious w.r.t. the leakage $|\mathbf{I}_1|$, $|\mathbf{I}_2|$, and RevealM . The running time is $O(NU \log(NU))$ where $N := |\mathbf{I}_1| + |\mathbf{I}_2|$ and $U = O(\frac{1}{\epsilon} \log \frac{1}{\delta})$.*

Proof. The obliviousness claim follows due to the fact that the oblivious send/receive algorithm is oblivious w.r.t. the lengths of the source and destination arrays; and the oblivious bin placement algorithm is oblivious w.r.t. the length of its input, as well as the bin capacities.

To see the running time, it suffices to show that the sum of all bin capacities is upper bounded by $O(NU)$. To see this, observe that there are at most N bin pairs, i.e., at most $2N$ bins. The total noise added to all bins is therefore upper bounded by $O(NU)$. Further, there are at most N real elements in all bins. Thus, the sum of the bin capacities is upper bounded by $O(NU)$. Note also that the length of L is $N = |\mathbf{I}_1| + |\mathbf{I}_2|$. Thus, the runtime claim follows due to the runtime of the oblivious send/receive algorithm and the oblivious bin placement algorithm (see Section 3.6). \square

4.3 Step 3: Bin-wise Cartesian Product

Once elements have been placed into bins like in Step 2, we can compute the Cartesian product of every pair of bins with the same join key, from \mathbf{I}_1 and \mathbf{I}_2 , respectively. This is done through the following procedure where $N := |\mathbf{I}_1| + |\mathbf{I}_2|$:

Step 3: Bin-wise Cartesian Product

For $i \in [N]$:

For every element (k, v) in bin $(1, i)$:

For every element (k', v') in bin $(2, i)$:

If $k = k' \in \mathcal{K}$ is a real join key, then add (k, v, v') to the output.

Else add the filler tuple (\perp, \perp, \perp) to the output.

Clearly, the access patterns of Step 3 are simulatable if we know the capacities of all bins, i.e., the `RevealM` array.

Fact 4.3. *Let `RevealM` be the same as L but removing the join key from each entry. Then, Step 3 is oblivious w.r.t. the leakage `RevealM`. Furthermore, if R is the true join-result length (i.e., when the result contains no filler entries), Step 3 runs in time at most $O(R + N \cdot U^2)$ where $U = \max \mathcal{G}(\epsilon, \delta) = O(\frac{1}{\epsilon} \log \frac{1}{\delta})$ and $N = |\mathbf{I}_1| + |\mathbf{I}_2|$.*

Proof. The obliviousness claim follows directly by the definition of the algorithm. To see the runtime, observe that the runtime of the algorithm is proportional to the length of the output. It therefore suffices to bound the length of the output, which can be obtained by summing up the occurrences of the following types of entries:

- *(Real, filler) entries.* In the output, each real element is matched with at most U filler elements. Thus, the total number of (real, filler) entries in the output is at most $O(NU)$.
- *(Filler, filler) entries.* There are at most N bin pairs, and each bin has at most U fillers. In total there are at most $O(NU)$ fillers. Each filler is matched with at most U fillers. Therefore, the total number of (filler, filler) entries in the output is at most $O(NU^2)$.
- *(Real, real) entries.* The number of (real, real) entries in the output is R .

\square

4.4 Step 4: Compaction of Result

Let Z be the output array from the above Step 3. Z may contain many filler elements of the form (\perp, \perp, \perp) . In one scan of Z , we can count the number of real elements in Z , and let R be this count. Now, we would like to compact Z to the length $R + \mathcal{G}(\epsilon, \delta, \Delta)$ where Δ is the maximum (noisy) count of any single join key (from either \mathbf{I}_1 or \mathbf{I}_2). This can be accomplished through the following algorithm.

Step 4: Compaction of Result

- (a) Let L'' be the set of noisy counts in L (pertaining to either \mathbf{I}_1 or \mathbf{I}_2). Let $\Delta := \max(L'')$.
- (b) Sample ξ from $\mathcal{G}(\epsilon, \delta, 2\Delta)$. Append a *fixed* $\widehat{U} := \max \mathcal{G}(\epsilon, \delta, 2\Delta)$ number of elements to the array Z : among the appended elements, ξ of them are *special filler* elements of the form $(\otimes, \otimes, \otimes)$ where $\otimes \notin \mathcal{K}$, and the remaining are filler elements of the form (\perp, \perp, \perp) .
- (c) Invoke a linear-time oblivious compaction algorithm (see Section 3.6) to compact the resulting array, moving all filler elements of the form (\perp, \perp, \perp) to the very end. We now remove the trailing part of the array containing filler elements of the form (\perp, \perp, \perp) . Note that the remaining array still has the special filler elements $(\otimes, \otimes, \otimes)$. Finally, in a linear scan of the outcome, rewrite all the special fillers $(\otimes, \otimes, \otimes)$ as (\perp, \perp, \perp) . Output the resulting array denoted \mathbf{O} .

Recall that the oblivious compaction algorithm leaks only the size of its input. Therefore, the above algorithm leaks only the length of the input Z , and the length of the output \mathbf{O} .

Fact 4.4. *Step 4 is oblivious w.r.t. the leakage $|Z|$ and $|\mathbf{O}|$. The running time of Step 4 is $O(|Z| + \frac{\Delta}{\epsilon} \log \frac{1}{\delta})$, and the output length is at most $O(R + \widehat{U})$.*

Proof. The obliviousness claim follows directly from the fact that the compaction algorithm is oblivious w.r.t. the input-length leakage, and the definition of the rest of the algorithm. The runtime is proportional to the length of the array output in Step 4(b), which is at most $|Z| + \widehat{U} = O(|Z| + \frac{\Delta}{\epsilon} \log \frac{1}{\delta})$. \square

Since Δ is a random variable upper bounded by $\mu_{\max} + U \leq N + U$ where μ_{\max} is the maximum multiplicity of any single element, and $|Z| = O(R + N \cdot U^2)$ due to Fact 4.3 the running time of Step 4 is at most $O(R + N \cdot U^2 + (N + U) \cdot \frac{1}{\epsilon} \log \frac{1}{\delta}) = O(R + N \cdot U^2 + (N + U)U) = O(R + N \cdot U^2)$ where $U = O(\frac{1}{\epsilon} \log \frac{1}{\delta})$. Further, the output length is at most $R + O(\widehat{U}) = R + O(\Delta \cdot \frac{1}{\epsilon} \log \frac{1}{\delta}) = R + O((\mu_{\max} + U) \cdot U)$.

Performance bounds of all steps. Summarizing all steps, the total runtime of the warmup algorithm is at most $O(R + NU^2 + NU \log(NU))$, and the result length is at most $R + O((\mu_{\max} + U) \cdot U)$.

4.5 Proof of Differential Obliviousness

By Lemma 3.10, it suffices to prove that the union of the leakages stated in Facts 4.1, 4.2, 4.3, 4.4 satisfy (ϵ', δ') -differential obliviousness. The union of the leakages in Facts 4.1, 4.2, 4.3 consists of `RevealM`, as well as $|\mathbf{I}_1|$ and $|\mathbf{I}_2|$. Fact 4.4 additionally leaks the output length besides the aforementioned leakage. Since $|\mathbf{I}_1|$ and $|\mathbf{I}_2|$ are publicly known given the input, we only care about the leakage `RevealM` and the output length.

Let $\text{Keys}(\mathbf{I}_1 \cup \mathbf{I}_2)$ be the union of distinct join keys appearing in either \mathbf{I}_1 or \mathbf{I}_2 , and let $\ell := |\mathbf{I}_1| + |\mathbf{I}_2| - |\text{Keys}(\mathbf{I}_1 \cup \mathbf{I}_2)|$. Recall that `RevealM` the multiset $\{(\widehat{n}_k^{(1)}, \widehat{n}_k^{(2)}) : k \in \text{Keys}(\mathbf{I}_1 \cup \mathbf{I}_2) \cup \mathcal{K}^*\}$ where $\mathcal{K}^* := \{\star_1, \dots, \star_\ell\}$. Each $\widehat{n}_k^{(b)}$ for $b \in \{1, 2\}$ is obtained by adding an independent $\mathcal{G}(\epsilon, \delta)$ noise to the true multiplicity of join key k in \mathbf{I}_b . Recall that `RevealM` is expressed in the implementation in a canonical form by lexicographical ordering of $(\widehat{n}_k^{(1)}, \widehat{n}_k^{(2)})$.

Lemma 4.5. *The leakage `RevealM` is $(2\epsilon, 2\delta)$ -differentially private where $N := |\mathbf{I}_1| + |\mathbf{I}_2|$.*

Proof. Let $\mathcal{I} := (\mathbf{I}_1, \mathbf{I}_2)$ and $\mathcal{I}' := (\mathbf{I}'_1, \mathbf{I}'_2)$ be two arbitrary neighboring inputs. We want to prove that $\text{RevealM}(\mathcal{I}) \sim_{(\epsilon, \delta)} \text{RevealM}(\mathcal{I}')$. We can consider two cases where $\text{Keys}(\mathcal{I} = (\mathbf{I}_1, \mathbf{I}_2)) := \text{Keys}(\mathbf{I}_1 \cup \mathbf{I}_2)$.

- Case 1: $\text{Keys}(\mathcal{I}) = \text{Keys}(\mathcal{I}')$. In this case, suppose we sort the join keys in $\text{Keys}(\mathcal{I}) \cup \mathcal{K}^*$ by lexicographical ordering, and let (k_1, \dots, k_N) denote the sorted sequence. Now, consider the ordered sequences

$$V := ((\hat{n}_{k_i}(\mathbf{I}_1), \hat{n}_{k_i}(\mathbf{I}_2)) : \text{for } i \in [N])$$

$$V' := ((\hat{n}_{k_i}(\mathbf{I}'_1), \hat{n}_{k_i}(\mathbf{I}'_2)) : \text{for } i \in [N])$$

where for $i \in [n]$ and $b \in \{1, 2\}$, $\hat{n}_{k_i}(\mathbf{I}_b)$ is the true count of k_i in \mathbf{I}_b plus a fresh sample of $\text{Geom}(\epsilon, \delta)$ noise.

There exists at most two pairs (b, i) where $b \in \{1, 2\}$ and $i \in [N]$ such that the true count of k_i in \mathbf{I}_b is not the same as the true count of k_i in \mathbf{I}'_b . For each such (b, i) , $\hat{n}_{k_i}(\mathbf{I}_b) \sim_{(\epsilon, \delta)} \hat{n}_{k_i}(\mathbf{I}'_b)$ by Fact 3.7. Now, by Facts 3.9, and 3.8, we have that $V \sim_{(2\epsilon, 2\delta)} V'$. Since $\text{RevealM}(\mathcal{I})$ and $\text{RevealM}(\mathcal{I}')$ can be obtained by applying a post-processing function to V and V' respectively reordering the array by the noisy counts it contains, we have $\text{RevealM}(\mathcal{I}) \sim_{(2\epsilon, 2\delta)} \text{RevealM}(\mathcal{I}')$ by Fact 3.8.

- Case 2: $\text{Keys}(\mathcal{I}) \neq \text{Keys}(\mathcal{I}')$. In this case, it must be that \mathcal{I} and \mathcal{I}' each has exactly one join key that the other does not have by the definition of neighboring — let $\mathbf{k} \in \mathcal{I}$ and $\mathbf{k}' \in \mathcal{I}'$ be the pair of join keys that the neighboring input does not have. Suppose that $\mathbf{k} \in \mathbf{I}_b$; in this case it must be that $\mathbf{k}' \in \mathbf{I}'_b$.

Now, let $\mathcal{K}^\bullet = (\text{Keys}(\mathcal{I}) \cup \mathcal{K}^*) \setminus \{\mathbf{k}\} = (\text{Keys}(\mathcal{I}') \cup \mathcal{K}^*) \setminus \{\mathbf{k}'\}$, and let

$$V := (\hat{n}_{\mathbf{k}}(\mathbf{I}_1), \hat{n}_{\mathbf{k}}(\mathbf{I}_2)), \quad \left(\text{for } \tilde{k} \in \mathcal{K}^\bullet \text{ in lexicographical order : } (\hat{n}_{\tilde{k}}(\mathbf{I}_1), \hat{n}_{\tilde{k}}(\mathbf{I}_2)) \right)$$

$$V' := (\hat{n}_{\mathbf{k}'}(\mathbf{I}'_1), \hat{n}_{\mathbf{k}'}(\mathbf{I}'_2)), \quad \left(\text{for } \tilde{k} \in \mathcal{K}^\bullet \text{ in lexicographical order : } (\hat{n}_{\tilde{k}}(\mathbf{I}'_1), \hat{n}_{\tilde{k}}(\mathbf{I}'_2)) \right)$$

It is not hard to see that V and V' are identically distributed.

Since $\text{RevealM}(\mathcal{I})$ and $\text{RevealM}(\mathcal{I}')$ can be obtained by applying a post-processing function to V and V' respectively reordering the array by the noisy counts it contains, we have $\text{RevealM}(\mathcal{I}) \sim_{(2\epsilon, 2\delta)} \text{RevealM}(\mathcal{I}')$ by Fact 3.8. □

Lemma 4.6. *Let $\Delta(\mathcal{I})$ be a random variable denoting the largest noisy count contained in $\text{RevealM}(\mathcal{I})$. Henceforth let $\text{supp}(\Delta(\mathcal{I}))$ denote the support of $\Delta(\mathcal{I})$ when the input is \mathcal{I} . Then, for any neighboring $\mathcal{I} := (\mathbf{I}_1, \mathbf{I}_2)$ and $\mathcal{I}' := (\mathbf{I}'_1, \mathbf{I}'_2)$, for any $\tilde{\Delta} \in \text{supp}(\Delta(\mathcal{I})) \cup \text{supp}(\Delta(\mathcal{I}'))$, it holds that*

$$\text{OutputLen}(\mathcal{I}, \tilde{\Delta}) \sim_{(\epsilon, \delta)} \text{OutputLen}(\mathcal{I}', \tilde{\Delta})$$

$\text{OutputLen}(\mathcal{I}, \tilde{\Delta})$ is the length of the true join result for input \mathcal{I} plus a noise sampled from the distribution $\mathcal{G}(\epsilon, \delta, 2\tilde{\Delta})$.

Proof. Observe that for neighboring \mathcal{I} and \mathcal{I}' , the true join result length cannot differ by more than any $2\tilde{\Delta}$ for any $\tilde{\Delta} \in \text{supp}(\Delta(\mathcal{I})) \cup \text{supp}(\Delta(\mathcal{I}'))$. Therefore, the lemma follows directly from Fact 3.7. □

Theorem 4.7 (Our warmup algorithm). *The warmup algorithm described in this section is $(3\epsilon, 3\delta + \text{negl}(N))$ -differentially oblivious where $\text{negl}(\cdot)$ is a negligible function¹¹. Furthermore, it completes*

¹¹If we use deterministic instantiation of the underlying oblivious sorting algorithm [AKS83], then the extra $\text{negl}(N)$ term will disappear.

in runtime

$$O\left(R + N \cdot \frac{1}{\epsilon} \log \frac{1}{\delta} \cdot \left(\log N + \frac{1}{\epsilon} \log \frac{1}{\delta}\right)\right)$$

and the result size is upper bounded by $R + O\left(\frac{1}{\epsilon}(\mu_{\max} + \frac{1}{\epsilon} \log \frac{1}{\delta}) \cdot \log \frac{1}{\delta}\right)$.

Proof. Since $\Delta(\mathcal{I})$ is fully determined by $\text{RevealM}(\mathcal{I})$, and by Lemma 4.5, $\text{RevealM}(\mathcal{I}) \sim_{(2\epsilon, 2\delta)} \text{RevealM}(\mathcal{I}')$ for any neighboring \mathcal{I} and \mathcal{I}' , it must be that $\Delta(\mathcal{I}) \sim_{(2\epsilon, 2\delta)} \Delta(\mathcal{I}')$ for any neighboring \mathcal{I} and \mathcal{I}' . It follows directly from Fact 3.9 and Lemma 4.6 that the combination of RevealM and OutputLen is $(3\epsilon, 3\delta)$ -differentially private.

By Facts 4.1, 4.2, 4.3, 4.4 the warmup algorithm's access patterns are $\text{negl}(N)$ -oblivious w.r.t. the leakage RevealM , the output length, and the lengths of the two input arrays. Now, the theorem follows from Lemma 3.10 and the fact that the combination of RevealM and OutputLen is $(3\epsilon, 3\delta)$ -differentially private.

The statement about the runtime can be attained by summing up the runtime stated in Facts 4.1, 4.2, 4.3, and 4.4 corresponding to the four steps. The claim about the result size follows directly from the definition of our noise distribution and the definition of Step 4. \square

4.6 A Cache Efficient Variant

As mentioned in Section 1, in Intel SGX-type application settings, the number of memory pages transmitted in and out of the trusted enclave is the dominant metric — this metric is also called the cache complexity in the standard algorithms literature. To instantiate a cache agnostic and cache efficient variant, we will

1. Instantiate the oblivious sorting algorithm with the randomized oblivious sorting algorithm by Ramachandran and Shi [RS20] which sorts n elements incurring $O(n \log n)$ runtime and $O((n/B) \log_{M/B}(n/B))$ cache complexity assuming the tall cache assumption that $M = \Omega(B^2)$ and further $M = \Omega(\log^{1+\epsilon} n)$ for an arbitrarily small constant $\epsilon \in (0, 1)$.
2. Instantiate the oblivious compaction algorithm with the scheme by Lin, Shi, and Xie [LSX19], which compacts an array of n elements incurring optimal, $O(n/B)$ cache complexity and $O(n \log \log n)$ runtime assuming that $M = \Omega(B^2)$ and $B \geq \log^{0.55} n$.

We now analyze the performance bounds of the cache-efficient version. Step 1 and Step 2 of the algorithm are dominated by $O(1)$ number of oblivious sorts on arrays of length at most $N \cdot U$ both in terms of runtime and cache complexity. Step 3, which computes the Cartesian product, incurs $O((R + N \cdot U^2)/B)$ cache complexity. Finally, Step 4 incurs $O((R + N \cdot U^2 + \widehat{U})/B)$ cache complexity where $U = \Theta(\frac{1}{\epsilon} \log \frac{1}{\delta})$ and $\widehat{U} = \Theta(\frac{\Delta}{\epsilon} \log \frac{1}{\delta}) = \Theta((\mu_{\max} + \frac{1}{\epsilon} \log \frac{1}{\delta}) \cdot \frac{1}{\epsilon} \log \frac{1}{\delta})$. Summarizing the above, the total cache complexity of all steps is upper bounded by

$$O\left(\frac{N}{B} \cdot \frac{1}{\epsilon} \log \frac{1}{\delta} \cdot \left(\log \frac{N}{B} + \frac{1}{\epsilon} \log \frac{1}{\delta}\right) + \frac{R}{B}\right)$$

Furthermore, the cache-efficient version incurs $O\left((R + N \cdot (\frac{1}{\epsilon} \log \frac{1}{\delta})^2) \cdot (\log \log N + \log \log(\frac{1}{\epsilon} \log \frac{1}{\delta}))\right)$ extra runtime due to the slightly non-optimal runtime of the oblivious compaction scheme by Lin et al. [LSX19].

5 Final Construction

Given our warmup idea, we now suggest an improved algorithm that further reduces the runtime to $O(R + N \log N)$. The intuitions of the final construction were explained earlier in Section 2.3, therefore, we jump directly into the detailed construction.

5.1 Step 1.1: Compute Bin Loads

Step 1.1 is the same as Step 1 of Section 4, except that now, in the output bin load array L , we also store the actual counts. In other words, every entry in L is of the form $(k, n_k^{(1)}, \hat{n}_k^{(1)}, n_k^{(2)}, \hat{n}_k^{(2)})$. To compute the new bin load array L , we can employ the same algorithm as in Section 4.1, and rather than overwriting the actual count with the noisy count, we shall remember both the actual count and noisy counts.

5.2 Step 1.2: Consolidate Sparse Join Keys

Given the load array L computed in Step 1.1, a join key k is said to be *sparse* iff its noise counts $\hat{n}_k^{(1)}$ and $\hat{n}_k^{(2)}$ are both upper bounded by $2U$, where $U = \max \mathcal{G}(\epsilon, \delta)$. A join key that is not sparse is said to be *dense*.

In Step 1.2, we will create an array `Bin` that stores the mapping from join keys to their destined bin-pair index. During this process, we map multiple sparse join keys to the same bin, such that each bin-pair will receive at least $\Theta(U)$ elements. In Step 1.2, we have not moved elements into their respective bins yet — the moving will be accomplished later in the binning step. We employ the following algorithm to compute the `BinMap` array; recall that $N := |\mathbf{I}_1| + |\mathbf{I}_2|$:

Step 1.2: Consolidate Sparse Join Keys

- (a) Define $L_s := \{(k, n_k^{(1)}, \hat{n}_k^{(1)}, n_k^{(2)}, \hat{n}_k^{(2)}) \in L : \hat{n}_k^{(1)} \leq 2U, \hat{n}_k^{(2)} \leq 2U\}$, i.e., L_s stores the multiplicities and noisy multiplicities of the sparse keys. Further, let $L_d := L \setminus L_s$. We shall assume that L_s and L_d are sorted according to the lexicographical order on $(\hat{n}_k^{(1)}, \hat{n}_k^{(2)})$ just like L . The arrays L_s and L_d can be constructed in a straightforward way in $O(N)$ time if we allow the access patterns to leak the noisy counts contained in L .
- (b) We will construct an array `BinMap` of length N that is initially empty. Elements of the form (k, i) will be appended to `BinMap` one by one where k is a join key appearing in L and i is some bin-pair index in the range $\{1, 2, \dots, m\}$ — we will show that $m = O(N/U)$.
- (c) For each $i \in \{1, 2, \dots, |L_d|\}$, if the i -th entry in L_d has the join key k , add the tuple (k, i) to the array `BinMap`. At this moment, `BinMap` stores the mapping from each *dense* join key to its destined bin-pair index.
- (d) Initially, let $j := |L_d| + 1$ to be the current bin-pair index. Let $c^{(1)}$ and $c^{(2)}$ be two counters that store the cumulative number of elements assigned to bin $(1, j)$ and bin $(2, j)$ respectively so far. Initially, both counters are 0.
- (e) Scan through each entry $(k, n_k^{(1)}, \hat{n}_k^{(1)}, n_k^{(2)}, \hat{n}_k^{(2)}) \in L_s$ in order:
 - if either $c^{(1)} + n_k^{(1)} > 4U$ or $c^{(2)} + n_k^{(2)} > 4U$, let $j := j + 1$ and reset both counters $c^{(1)}$ and $c^{(2)}$ to be 0;
 - if at least one of $n_k^{(1)}$ and $n_k^{(2)}$ is non-zero, append the tuple (k, j) `BinMap`, and let

$$c^{(b)} := c^{(b)} + n_k^{(b)} \text{ for } b \in \{1, 2\}; \text{ otherwise append the tuple } (k, \perp) \text{ to BinMap.}$$

(f) Finally, output the array **BinMap** that maps each join key to its bin-pair index, and also output **DenseCapacities** := $\{(\hat{n}_k^{(1)}, \hat{n}_k^{(2)}) : (k, n_k^{(1)}, \hat{n}_k^{(1)}, n_k^{(2)}, \hat{n}_k^{(2)}) \in L_d\}$.

Henceforth, a bin-pair designated for a single dense join key is called a *D-bin-pair* and the bins are called *D-bins*. A bin-pair designated for possibly multiple sparse join keys is called an *S-bin-pair*, and the bins are called *S-bins*. We first need to verify that the above algorithm indeed consumes at most $O(N/U)$ bin pairs, as stated in the following Fact 5.2. In this way, in the subsequent binning step, we will just set the total number of bins to an a-priori fixed upper bound CN/U for a sufficiently large constant C . It is possible that some S-bins have no elements mapped to them by the **BinMap** array.

Fact 5.1. *At the end of the above algorithm, the bin-pair index j is upper bounded by $O(N/U)$.*

Proof. By the definition of a dense join key and also given that our noise is bounded in the range $[0, U]$, any dense join key must have a total multiplicity of U in $\mathbf{I}_1 \cup \mathbf{I}_2$. Thus, there can be at most $O(N/U)$ D-bin-pairs.

It remains to show that there are at most $O(N/U)$ S-bin-pairs. By the definition of a sparse join key and also given that our noise is bounded in the range $[0, U]$, any sparse join key can have at most $2U$ multiplicities in either \mathbf{I}_1 or \mathbf{I}_2 . We open up a new bin by incrementing j if adding a new join key k to the present bin j will exceed the bin's capacity $4U$. Therefore, before we open up a new bin, the current bin's capacity must be at least $2U$. This means that the number of S-bins is upper bounded by $N/2U$. \square

Fact 5.2. *The sparse join key consolidation step is oblivious w.r.t. the leakage N and **RevealM**. The running time is $O(N)$.*

Proof. The proof is straightforward and can be checked by going through the algorithm line by line. \square

5.3 Step 2: Binning

The binning step can be done similarly as before, but now using the following plan for the binning:

1. The mapping from join keys to their respective bins is specified by the **BinMap** array: specifically, suppose $(k, i) \in \mathbf{BinMap}$, this means that elements with the join key k from \mathbf{I}_b should go into bin (b, i) for $b \in \{1, 2\}$.
2. There are a total of CN/U bins where C is a sufficiently large constant. The first $|L_d|$ of them are D-bin-pairs whose capacities are specified by **DenseCapacities**, whereas the remaining are S-bin-pairs, and each S-bin has capacity exactly $4U$.

Note that because we padded the number of S-bins to an a-priori fixed upper bound, it is possible that according to **BinMap**, some S-bins are not associated with any join keys.

To perform the binning, we first rely on oblivious send-recv such that each real element in \mathbf{I}_1 and \mathbf{I}_2 receives its bin index, and then we rely on oblivious bin placement to place the real elements into bins.

Fact 5.3. *The binning step is oblivious w.r.t. the leakage $|\mathbf{I}_1|$, $|\mathbf{I}_2|$ and **RevealM**. The running time is $O(N \log N)$.*

Proof. For the claim about access patterns, observe that the access patterns of oblivious send-receive can be simulated given $|\text{BinMap}|$, $|\mathbf{I}_1|$, and $|\mathbf{I}_2|$. Further, the access patterns of oblivious bin placement can be simulated given $|\text{BinMap}|$, $|\mathbf{I}_1|$, and $|\mathbf{I}_2|$, and the capacities of all bins. Recall also that the capacities of bins for sparse join keys are $4U$ and the capacities of bins for dense join keys are determined by `RevealM`.

Due to the proof of Fact 5.2, fillers occupy at most a constant fraction of each bin pair. Therefore, the total sum of capacities over all bins is $O(N)$. The claim on running time now follows from the running time of oblivious send-receive and oblivious bin placement (see Section 3.6). \square

5.4 Remaining Steps

Perform the same Steps 3-4 as in Section 4. During the Cartesian product step, our earlier warmup algorithm will not try to join two real elements with different join keys; but now, since each bin can contain multiple join keys, this can happen. Fortunately, our previous Cartesian product algorithm still works, since by the algorithm description, if we try to join two real elements with different join keys, a filler element will be written to the output.

Fact 5.4. *The bin-wise Cartesian product step is oblivious w.r.t. the leakage `RevealM` and N ; further, it has a running time of $O(R + NU)$, where $U = O(\frac{1}{\epsilon} \log \frac{1}{\delta})$.*

Proof. When we compute the bin-wise Cartesian product, each real element in \mathbf{I}_b where $b \in \{1, 2\}$ is paired with at most $O(U)$ filler elements or elements with different join keys from \mathbf{I}_{3-b} — this introduces at most NU fillers into the bin-wise Cartesian product result. By the definition of the algorithm, each bin contains at most $O(U)$ filler elements. Moreover, each filler element in \mathbf{I}_b can be paired with at most $O(U)$ filler elements from \mathbf{I}_{3-b} — by Fact 5.2, since the number of bins is $O(N/U)$, the filler-filler pairs introduces another $O(NU)$ elements into the bin-wise Cartesian product result. Thus, the total number of filler elements in the result is upper bounded by $O(NU)$.

The running time of the bin-wise Cartesian product step is upper bounded by some constant times the result length, which is at most $O(R + NU)$. \square

Theorem 5.5 (Our DO join algorithm). *The final join algorithm described in this section is $(3\epsilon, 3\delta + \text{negl}(N))$ -differentially oblivious where $\text{negl}(\cdot)$ is a negligible function, runs in time $O(R + N(\log N + \frac{1}{\epsilon} \log \frac{1}{\delta}))$ and produces a result whose length is at most $R + O(\frac{1}{\epsilon}(\mu_{\max} + \frac{1}{\epsilon} \log \frac{1}{\delta}) \cdot \log \frac{1}{\delta})$.*

Proof. The running time follows from Facts 5.2, 5.3 and 5.4. The analysis of the result length is the same as in Theorem 4.7. Since the algorithm is $\text{negl}(N)$ -oblivious w.r.t. the leakage $|\mathbf{I}_1|$, $|\mathbf{I}_1|$, `RevealM` and the final output length, differential obliviousness follows the exact argument in Theorem 4.7. \square

5.5 A Cache Efficient Variant

To obtain a cache agnostic, cache efficient variant, we can switch the oblivious sort and oblivious compaction building blocks to cache agnostic, cache efficient variants just like in Section 4.6.

Steps 1.1, 1.2, and Step 2 are dominated, both in terms of runtime and cache complexity, by a constant number of oblivious sorts on arrays of size at most $O(N)$. Step 3, i.e, Cartesian product step, now produces an array of length at most $O(R + NU)$, and thus its cache complexity is $O((R + NU)/B)$. Step 4, compaction of result, has cache complexity upper bounded $O((R + NU + \widehat{U})/B)$.

Therefore, the cache complexity of all steps is upper bounded by

$$O\left(\frac{R}{B} + \frac{N}{B} \cdot \left(\log_{\frac{M}{B}} \frac{N}{B} + \frac{1}{\epsilon} \log \frac{1}{\delta}\right) + \frac{1}{B} \left(\frac{1}{\epsilon} \log \frac{1}{\delta}\right)^2\right)$$

This cache efficient variant incurs $O\left((R + N \cdot \frac{1}{\epsilon} \log \frac{1}{\delta}) \cdot (\log \log N + \log \log(\frac{1}{\epsilon} \log \frac{1}{\delta}))\right)$ extra runtime due to the slightly non-optimal runtime of the oblivious compaction scheme by Lin et al. [LSX19].

Summarizing the above analysis, we obtain the following corollary:

Corollary 5.6 (Our DO join algorithm: cache complexity). *There is an (ϵ, δ) -DO database join algorithm that incurs cache complexity upper bounded by $\frac{1}{B} \cdot O\left(N \left(\log_{\frac{M}{B}} \frac{N}{B} + \frac{1}{\epsilon} \log \frac{1}{\delta}\right) + R + \left(\frac{1}{\epsilon} \log \frac{1}{\delta}\right)^2\right)$, assuming the tall cache assumption $M = \Omega(B^2)$ and the wide block assumption $B = \Omega(\log^{0.55} N)$.*

6 Lower Bounds

6.1 The Complexity of Fully Oblivious Database Join

It is not hard to show that the complexity of fully oblivious database join is $\Theta(|\mathbf{I}_1| \cdot |\mathbf{I}_2|)$ even when the actual join answer (without fillers) can be much shorter. Furthermore, this is both the upper and lower bound. Intuitively, this is because a fully oblivious join must hide the multiplicity of each join key within each input array. More formally, we prove the following theorem:

Theorem 6.1. *Any fully oblivious database join algorithm must, on any input instance, incur $\Omega(|\mathbf{I}_1| \cdot |\mathbf{I}_2|)$ runtime except with negligible probability. Furthermore, there exists a fully oblivious database join algorithm that incurs $O(|\mathbf{I}_1| \cdot |\mathbf{I}_2|)$ runtime.*

Proof. With a fully oblivious database join, the output's length must be statistically indistinguishable for any two inputs $(\mathbf{I}_1, \mathbf{I}_2)$ and $(\mathbf{J}_1, \mathbf{J}_2)$ of equal length. Now, it could be the case that all of \mathbf{I}_1 and \mathbf{I}_2 have the same join key k . In this case, for the output to be correct and not lose any element in the Cartesian product, its length must be at least $|\mathbf{I}_1| \cdot |\mathbf{I}_2|$. Therefore, if the algorithm is oblivious, its output length should be at least $|\mathbf{I}_1| \cdot |\mathbf{I}_2|$ on any input except with negligible probability.

A naïve algorithm that can achieve fully oblivious database join with runtime $O(|\mathbf{I}_1| \cdot |\mathbf{I}_2|)$ was described earlier in Section 3.5. \square

For the special case when $|\mathbf{I}_1| = |\mathbf{I}_2| = N$, the running time would be $\Theta(N^2)$. Note that the length of the true answer (i.e., without any fillers), denoted R , could be much smaller than N^2 . In this case, an insecure algorithm can compute the join in $O(R + N)$ time (see Section 3.5).

6.2 Lower Bound for Differentially Oblivious Database Join

6.2.1 Lower Bound on Result Size

Fact 6.2. *Let $\epsilon > 0$ and $\delta \in (0, 1)$ such that $\epsilon \geq \delta$ and $n = \Omega(\frac{1}{\epsilon} \log \frac{\epsilon}{\delta})$. Suppose $\{F_i : i \in [n]\}$ is a collection of random variables such that for all $i, j \in [n]$, $|i - j| \leq 1$ implies that $F_i \sim_{(\epsilon, \delta)} F_j$. Suppose E and B satisfy that for all $i \in [n]$, $|F_i - i \cdot B| \leq E$ holds with probability at least $1 - \gamma$. Then, $E = \Omega(\frac{B}{\epsilon} \min\{\log \frac{\epsilon}{\delta}, \log \frac{1}{\gamma}\})$.*

Proof. A standard calculation gives that $|i - j| \leq k$ implies that $F_i \sim_{(k\epsilon, e^{k\epsilon}, \frac{\delta}{\epsilon})} F_j$. Let S_0 be the support of F_0 . Observe that $\Pr[F_0 \in S_0] = 1$, but the hypothesis implies that for $K = \lceil \frac{2E}{B} \rceil$, $\Pr[F_K \in S_0] \leq \beta$. However, since $F_0 \sim_{(K\epsilon, e^{K\epsilon}, \frac{\delta}{\epsilon})} F_K$, we must have the following, as required:

$$1 \leq e^{K\epsilon} \cdot \gamma + e^{K\epsilon} \cdot \frac{\delta}{\epsilon}.$$

□

The theorem below says that any (ϵ, δ) -differentially oblivious algorithm must have an input such that the algorithm produces an output of length at least $R + \Omega(\mu_{\max} \cdot \frac{1}{\epsilon} \cdot \min\{\log \frac{\epsilon}{\delta}, \log \frac{1}{\gamma}\})$ with probability at least γ . For sufficiently small γ , the lower bound becomes $R + \Omega(\mu_{\max} \cdot \frac{1}{\epsilon} \cdot \log \frac{\epsilon}{\delta})$. For the special case when $\mu_{\max} = \Theta(N)$, $\epsilon = \Theta(1)$, and $\delta < 1/N$, we thus have an $R + \Omega(N \log N)$ lower bound on the result length.

Theorem 6.3 (Lower bound on the result size of any DO join algorithm). *Suppose μ_{\max} is the maximum multiplicity of a join key in either \mathbf{I}_1 or \mathbf{I}_2 . Let $\epsilon > 0$ and $\delta \in (0, 1)$ such that $\epsilon \geq \delta$ and $|\mathbf{I}_1| + |\mathbf{I}_2| \geq \Omega(\frac{1}{\epsilon} \log \frac{\epsilon}{\delta})$. Moreover, suppose an (ϵ, δ) -differentially oblivious join algorithm, for any input, outputs a result of size at most $R + E$ with probability at least $1 - \gamma$, where R is the true join result size. Then, $E \geq \Omega(\frac{\mu_{\max}}{\epsilon} \cdot \min\{\log \frac{\epsilon}{\delta}, \log \frac{1}{\gamma}\})$.*

In other words, there must exist an input for which the algorithm produces an output of length at least $R + \Omega(\frac{\mu_{\max}}{\epsilon} \cdot \min\{\log \frac{\epsilon}{\delta}, \log \frac{1}{\gamma}\})$ with probability at least γ .

Proof. Fix sufficiently large $n = \Omega(\frac{1}{\epsilon} \log \frac{\epsilon}{\delta})$. For $i \in [n]$, define an input $\mathcal{I}^{(i)}$ such that there are only two join keys k_1 and k_2 : in array \mathbf{I}_1 , k_1 has multiplicity i and k_2 has multiplicity $n - i$; in array \mathbf{I}_2 , k_1 has multiplicity $B := \mu_{\max}$, and k_2 has multiplicity 0.

Hence, the correct output length for $\mathcal{I}^{(i)}$ is $i \cdot B$. Define F_i to be the length of the output on $\mathcal{I}^{(i)}$ given by the differentially oblivious mechanism. Fact 6.2 gives the required lower bound. □

6.2.2 Lower Bound on Runtime

We prove a lower bound on the runtime of any (ϵ, δ) -DO join algorithm that follows the so-called *indivisible model*. In the indivisible model, the algorithm is allowed to 1) make a copy of an element's payload, and 2) move the payload string around in memory; but it is not allowed to perform encoding or computation on the payload string. Most natural algorithms (not just for join, but also for other computational tasks such as sorting) fit the indivisible model.

Our formal theorem statement below gives a generalized parametrization, but observe that for the typical case when $\epsilon = \Theta(1)$ and $\delta < 1/N$, we obtain a runtime lower bound of $\Omega(R + N \log \log N + \mu_{\max} \log N)$. Another intriguing and initially counter-intuitive implication of the following theorem is that $\Omega(R + N \log N)$ runtime is necessary for any $(\epsilon, 0)$ -DO join algorithm even when ϵ can be arbitrarily large — even though when ϵ is very large, there is almost no privacy left. The reason for this is because if $\delta = 0$, then any access pattern that is encountered with non-zero probability for one input $(\mathbf{I}_1, \mathbf{I}_2)$ must be encountered with non-zero probability for *any* other input $(\mathbf{I}'_1, \mathbf{I}'_2)$ as long as $|\mathbf{I}_1| = |\mathbf{I}'_1|$ and $|\mathbf{I}_2| = |\mathbf{I}'_2|$. In other words, there must exist a feasible plan for routing the elements' payload strings along the access pattern graph such that the outcome is a correct join result. This imposes a large lower bound on the size of the access pattern graph which translates to a lower bound on the algorithm's runtime.

Theorem 6.4 (Lower bound on the runtime of any DO join algorithm). *Let N be the total input length and let $0 \leq s \leq \sqrt{N}$ be an integer. Suppose $\epsilon > \delta > 0$, $0 \leq \delta \leq \beta \cdot \frac{\epsilon}{s} \cdot e^{-2cs}$ and $N \geq \Omega(\frac{1}{\epsilon} \log \frac{\epsilon}{\delta})$.*

Then, any join algorithm in the indivisible model that is (ϵ, δ) -differentially oblivious must have some input of total length N such that the algorithm incurs at least $\Omega(R + N \log s + \mu_{\max} \cdot \frac{1}{\epsilon} \cdot \min\{\log \frac{\epsilon}{\delta}, \log \frac{1}{1-\beta}\})$ runtime with probability at least $1 - 2\beta$, where R denotes the true result length and μ_{\max} is the maximum multiplicity of any join key in either input array.

Proof. Let Join be an arbitrary (ϵ, δ) -DO join algorithm. Due to Theorem 6.3, there must be an input such that it takes Join at least $\Omega(R + \mu_{\max} \cdot \frac{1}{\epsilon} \cdot \log \frac{\epsilon}{\delta} \min\{\log \frac{\epsilon}{\delta}, \log \frac{1}{1-\beta}\})$ time to write down the output with probability at least $1 - \beta$ — observe that Theorem 6.3 applies to not just the result length, but also the number of memory locations in the output array that the algorithm must visit.

Below we prove that there is some input which would cost Join at least $\Omega(N \log s)$ runtime with probability $1 - \beta$. If we can prove this, this would imply the stated theorem. To prove the above claim, our idea is to show a privacy-preserving reduction from sorting to join. Specifically, we will show that if there is an (ϵ, δ) -DO join algorithm Join in the indivisible model that guarantees $T(N, \epsilon, \delta)$ runtime on any input in which the two input arrays have equal length $N/2$, then there is an (ϵ, δ) -DO sorting algorithm that sorts N elements in time $C \cdot T(N, \epsilon, \delta)$ for some suitable constant $C \geq 1$. If we can show this, then the theorem follows directly due to a lower bound for DO sorting shown by Chan et al. [CCMS19] (Theorem 4.7 and Corollary 4.8 in their paper).

The reduction works as follows. Henceforth, we may assume that in the sorting problem, the N join keys in the input array \mathbf{I} are all distinct and take value from the domain $[N]$, since the DO sorting lower bound by Chan et al. [CCMS19] holds for this special case of sorting. To sort the input array \mathbf{I} , we will leverage the Join algorithm as follows. Let $\mathbf{I}_1 := \mathbf{I}$ and let \mathbf{I}_2 be an array in which the elements have join keys $1, 2, \dots, N$ respectively, and the payload is a canonical string \star that indicates that the element comes from \mathbf{I}_2 .

Now, we run Join on $(\mathbf{I}_1, \mathbf{I}_2)$. Suppose that each step t of the algorithm performs reads the memory address raddr_t and writes to the memory address waddr_t . We can augment the Join algorithm such that in each step t of the algorithm, it appends the operation performed in this step to an array denoted ops . Specifically, the recorded operation in step t is a tuple of the form $(\text{raddr}_t, \text{waddr}_t, r, w)$. Besides the raddr_t and waddr_t fields defined earlier, if $r \neq \perp$ it means that some element's payload string is moved from raddr_t into the r -th register; and if $w \neq \perp$ it means that the payload string stored in the w -th register is moved into memory location waddr_t .

After the augmented Join algorithm completes, the two elements with the same join key coming from \mathbf{I}_1 and \mathbf{I}_2 respectively are placed together in the output array. By linearly scanning through the output array, we can exchange the two elements with the same join key coming from \mathbf{I}_1 and \mathbf{I}_2 respectively.

At this moment, we replay the reverse of the operations stored in ops in reverse order. If a stored entry is of the form $(\text{raddr}_t, \text{waddr}_t, r, w)$, we will now perform the reverse operation, that is,

- If $r \neq \perp$, write the element payload stored in register r to memory location raddr_t . Else, if $r = \perp$, make a fake write to memory location raddr_t .
- If $w \neq \perp$, read the element payload stored memory location waddr_t to register w . Else, if $w = \perp$, perform a fake read of memory location waddr_t .

After replaying the reverse of the operations in reverse order, the elements' payload strings originally stored in \mathbf{I}_1 will be moved to \mathbf{I}_2 and the outcome is in sorted order.

Since the access pattern revealed by the above sorting algorithm is simply the access patterns of Join plus its reverse, if Join is (ϵ, δ) -DO, so is the above sorting algorithm. This completes our proof. □

Acknowledgments

We gratefully acknowledge helpful discussions and insightful feedback with Zhao Song and Lianke Qin. T-H. Hubert Chan is partially supported by the Hong Kong RGC under the grants 17200418 and 17201220. Elaine Shi is partially funded by NSF under award numbers 1601879 and 2128519, a faculty award from JP Morgan, an ONR YIP award, and a Packard Fellowship.

References

- [ABD⁺07] Lars Arge, Michael A. Bender, Erik D. Demaine, Bryan Holland-Minkley, and J. Ian Munro. An optimal cache-oblivious priority queue and its application to graph algorithms. *SIAM Journal on Computing*, 36(6):1672–1695, 2007.
- [ABE⁺13] Arvind Arasu, Spyros Blanas, Ken Eguro, Raghav Kaushik, Donald Kossmann, Ravishankar Ramamurthy, and Ramarathnam Venkatesan. Orthogonal security with cipherbase. In *CIDR*, 2013.
- [ACN⁺20] Gilad Asharov, T.-H. Hubert Chan, Kartik Nayak, Rafael Pass, Ling Ren, and Elaine Shi. Bucket oblivious sort: An extremely simple oblivious sort. In Martin Farach-Colton and Inge Li Gørtz, editors, *3rd Symposium on Simplicity in Algorithms, SOSA@SODA*, pages 8–14. SIAM, 2020.
- [AFS19] Prabhanjan Ananth, Xiong Fan, and Elaine Shi. Towards attribute-based encryption for RAMs from LWE: sub-linear decryption, and more. In Steven D. Galbraith and Shiho Moriai, editors, *Advances in Cryptology - ASIACRYPT 2019 - 25th International Conference on the Theory and Application of Cryptology and Information Security, Kobe, Japan, December 8-12, 2019, Proceedings, Part I*, volume 11921 of *Lecture Notes in Computer Science*, pages 112–141. Springer, 2019.
- [AGM08] Albert Atserias, Martin Grohe, and Dániel Marx. Size bounds and query plans for relational joins. In *FOCS*, page 739–748, USA, 2008. IEEE Computer Society.
- [AHKM19] Archita Agarwal, Maurice Herlihy, Seny Kamara, and Tarik Moataz. Encrypted databases for differential privacy. *Proc. Priv. Enhancing Technol.*, 2019(3):170–190, 2019.
- [AHV95] Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [AK14] Arvind Arasu and Raghav Kaushik. Oblivious query processing. In *Proc. 17th International Conference on Database Theory (ICDT), Athens, Greece, March 24-28, 2014*, pages 26–37. OpenProceedings.org, 2014.
- [AKL⁺20a] Gilad Asharov, Ilan Komargodski, Wei-Kai Lin, Kartik Nayak, Enoch Peserico, and Elaine Shi. OptORAMA: Optimal Oblivious RAM. In *Advances in Cryptology - EUROCRYPT 2020*, 2020. To appear. See also: <https://eprint.iacr.org/2018/892>.
- [AKL⁺20b] Gilad Asharov, Ilan Komargodski, Wei-Kai Lin, Enoch Peserico, and Elaine Shi. Oblivious parallel tight compaction. In *Information-Theoretic Cryptography (ITC)*, 2020.

- [AKS83] M. Ajtai, J. Komlós, and E. Szemerédi. An $O(n \log n)$ sorting network. In *STOC*, 1983.
- [ANS09] Yuriy Arbitman, Moni Naor, and Gil Segev. De-amortized cuckoo hashing: Provable worst-case performance and experimental results. In *ICALP*, 2009.
- [AV88] Alok Aggarwal and S. Vitter, Jeffrey. The Input/Output Complexity of Sorting and Related Problems. *Commun. ACM*, 31(9):1116–1127, September 1988.
- [BCP15a] Elette Boyle, Kai-Min Chung, and Rafael Pass. Large-scale secure computation: Multi-party computation for (parallel) RAM programs. In *Advances in Cryptology - CRYPTO 2015 - 35th Annual Cryptology Conference, Santa Barbara, CA, USA, August 16-20, 2015, Proceedings, Part II*, volume 9216 of *Lecture Notes in Computer Science*, pages 742–762. Springer, 2015.
- [BCP15b] Elette Boyle, Kai-Min Chung, and Rafael Pass. Oblivious parallel ram. In *Theory of Cryptography Conference (TCC)*, 2015.
- [BDB16] Jeremiah Blocki, Anupam Datta, and Joseph Bonneau. Differentially private password frequency lists. In *NDSS*, 2016.
- [BKS17] Paul Beame, Paraschos Koutris, and Dan Suciu. Communication steps for parallel query processing. *J. ACM*, 64(6):40:1–40:58, 2017.
- [BNS16] Mark Bun, Kobbi Nissim, and Uri Stemmer. Simultaneous private learning of multiple concepts. In *Proceedings of the 2016 ACM Conference on Innovations in Theoretical Computer Science (ITCS)*, pages 369–380, 2016.
- [BNZ19] Amos Beimel, Kobbi Nissim, and Mohammad Zaheri. Exploring differential obliviousness. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2019, September 20-22, 2019, Massachusetts Institute of Technology, Cambridge, MA, USA*, volume 145 of *LIPIcs*, pages 65:1–65:20, 2019.
- [BPWZ14] Andreas Björklund, Rasmus Pagh, Virginia Vassilevska Williams, and Uri Zwick. Listing Triangles. In *Proceedings of the 41st International Colloquium on Automata, Languages, and Programming (ICALP)*, volume 8572 of *Lecture Notes in Computer Science*, pages 223–234. Springer International Publishing, 2014.
- [BS14] Sumeet Bajaj and Radu Sion. Trusteddb: A trusted hardware-based database with privacy and data confidentiality. *IEEE Trans. on Knowl. and Data Eng.*, 26(3):752–765, March 2014.
- [BV18] Victor Balcer and Salil P. Vadhan. Differential privacy on finite computers. In Anna R. Karlin, editor, *9th Innovations in Theoretical Computer Science Conference (ITCS)*, volume 94, pages 43:1–43:21. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018.
- [CBC⁺18] Natacha Crooks, Matthew Burke, Ethan Cecchetti, Sitar Harel, Rachit Agarwal, and Lorenzo Alvisi. Obladi: Oblivious serializable transactions in the cloud. In *OSDI*, page 727–743, USA, 2018. USENIX Association.

- [CBS15] Shumo Chu, Magdalena Balazinska, and Dan Suciu. From theory to practice: Efficient join query evaluation in a parallel database system. In *SIGMOD Conference*, pages 63–78. ACM, 2015.
- [CCC⁺16] Yu-Chi Chen, Sherman S. M. Chow, Kai-Min Chung, Russell W. F. Lai, Wei-Kai Lin, and Hong-Sheng Zhou. Cryptography for parallel RAM from indistinguishability obfuscation. In *Proceedings of the 2016 ACM Conference on Innovations in Theoretical Computer Science, Cambridge, MA, USA, January 14-16, 2016*, pages 179–190. ACM, 2016.
- [CCHR16] Ran Canetti, Yilei Chen, Justin Holmgren, and Mariana Raykova. Adaptive succinct garbled RAM or: How to delegate your database. In *Theory of Cryptography - 14th International Conference, TCC 2016-B, Beijing, China, October 31 - November 3, 2016, Proceedings, Part II*, volume 9986 of *Lecture Notes in Computer Science*, pages 61–90, 2016.
- [CCMS19] T.-H. Hubert Chan, Kai-Min Chung, Bruce M. Maggs, and Elaine Shi. Foundations of differentially oblivious algorithms. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, page 2448–2467, USA, 2019. Society for Industrial and Applied Mathematics.
- [CCS17] Hubert Chan, Kai-Min Chung, and Elaine Shi. On the depth of oblivious parallel oram. manuscript, 2017.
- [CFK⁺15] Marek Cygan, Fedor V. Fomin, Lukasz Kowalik, Daniel Lokshtanov, Daniel Marx, Marcin Pilipczuk, Michal Pilipczuk, and Saket Saurabh. *Parameterized Algorithms*. Springer Publishing Company, Incorporated, 1st edition, 2015.
- [CGKO06] Reza Curtmola, Juan A. Garay, Seny Kamara, and Rafail Ostrovsky. Searchable symmetric encryption: improved definitions and efficient constructions. In *ACM Conference on Computer and Communications Security*, pages 79–88, 2006.
- [CGPR15] David Cash, Paul Grubbs, Jason Perry, and Thomas Ristenpart. Leakage-abuse attacks against searchable encryption. In *ACM CCS*, page 668–679, 2015.
- [CJJ⁺13] David Cash, Stanislaw Jarecki, Charanjit S. Jutla, Hugo Krawczyk, Marcel-Catalin Rosu, and Michael Steiner. Highly-scalable searchable symmetric encryption with support for boolean queries. In *CRYPTO*, 2013.
- [CJJ⁺14] David Cash, Joseph Jaeger, Stanislaw Jarecki, Charanjit S. Jutla, Hugo Krawczyk, Marcel-Catalin Rosu, and Michael Steiner. Dynamic searchable encryption in very-large databases: Data structures and implementation. In *21st Annual Network and Distributed System Security Symposium, NDSS 2014, San Diego, California, USA, February 23-26, 2014*. The Internet Society, 2014.
- [CKX06] Jianer Chen, Iyad A. Kanj, and Ge Xia. Improved parameterized upper bounds for vertex cover. In *In MFCS*, 2006.
- [CS17] T.-H. Hubert Chan and Elaine Shi. Circuit OPRAM: unifying statistically and computationally secure ORAMs and OPRAMs. In *Theory of Cryptography Conference, (TCC)*, 2017.

- [CSS11] T.-H. Hubert Chan, Elaine Shi, and Dawn Song. Private and continual release of statistics. *TISSEC*, 14(3):26, 2011.
- [Dem02] Erik D. Demaine. Cache-oblivious algorithms and data structures. In *Lecture Notes from the EEF Summer School on Massive Data Sets*. BRICS, University of Aarhus, Denmark, June 27–July 1 2002.
- [DMNS06] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, 2006.
- [DR14] Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4):211–407, 2014.
- [EZ19] Saba Eskandarian and Matei Zaharia. Oblidb: Oblivious query processing for secure databases. *Proc. VLDB Endow.*, 13(2):169–183, 2019.
- [FG06] J. Flum and M. Grohe. *Parameterized Complexity Theory (Texts in Theoretical Computer Science. An EATCS Series)*. Springer-Verlag, Berlin, Heidelberg, 2006.
- [FLPR99] Matteo Frigo, Charles E Leiserson, Harald Prokop, and Sridhar Ramachandran. Cache-oblivious algorithms. In *Foundations of Computer Science, 1999. 40th Annual Symposium on*, pages 285–297. IEEE, 1999.
- [GHL⁺14] Craig Gentry, Shai Halevi, Steve Lu, Rafail Ostrovsky, Mariana Raykova, and Daniel Wichs. Garbled ram revisited. In *EUROCRYPT*, pages 405–422, 2014.
- [GHRW14] Craig Gentry, Shai Halevi, Mariana Raykova, and Daniel Wichs. Outsourcing private ram computation. In *STOC*, 2014.
- [GKK⁺12] S. Dov Gordon, Jonathan Katz, Vladimir Kolesnikov, Fernando Krell, Tal Malkin, Mariana Raykova, and Yevgeniy Vahlis. Secure two-party computation in sublinear (amortized) time. In *CCS*, 2012.
- [GKP⁺13] Shafi Goldwasser, Yael Tauman Kalai, Raluca A. Popa, Vinod Vaikuntanathan, and Nikolai Zeldovich. How to run turing machines on encrypted data. In Ran Canetti and Juan A. Garay, editors, *Advances in Cryptology - CRYPTO 2013 - 33rd Annual Cryptology Conference, Santa Barbara, CA, USA, August 18-22, 2013. Proceedings, Part II*, volume 8043 of *Lecture Notes in Computer Science*, pages 536–553. Springer, 2013.
- [GLMP19] Paul Grubbs, Marie-Sarah Lacharité, Brice Minaud, and Kenneth G. Paterson. Learning to reconstruct: Statistical learning theory and encrypted database attacks. In *2019 IEEE Symposium on Security and Privacy, SP 2019, San Francisco, CA, USA, May 19-23, 2019*, pages 1067–1083. IEEE, 2019.
- [GMN⁺16] Paul Grubbs, Richard McPherson, Muhammad Naveed, Thomas Ristenpart, and Vitaly Shmatikov. Breaking web applications built on top of encrypted data. In *ACM CCS*, page 1353–1364, 2016.
- [GO96] Oded Goldreich and Rafail Ostrovsky. Software protection and simulation on oblivious RAMs. *J. ACM*, 1996.

- [Gol87] O. Goldreich. Towards a theory of software protection and simulation by oblivious RAMs. In *STOC*, 1987.
- [GRU12] Anupam Gupta, Aaron Roth, and Jonathan Ullman. Iterative constructions and private data release. In *TCC*, volume 7194, pages 339–356, 2012.
- [GUW09] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. *Database systems - the complete book (2. ed.)*. Pearson Education, 2009.
- [HR10] Moritz Hardt and Guy N. Rothblum. A multiplicative weights mechanism for privacy-preserving data analysis. In *FOCS*, pages 61–70, 2010.
- [HY19] Xiao Hu and Ke Yi. Instance and output optimal parallel algorithms for acyclic joins. In *Proceedings of the 38th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS)*, page 450–463, 2019.
- [IKK12] Mohammad Islam, Mehmet Kuzu, and Murat Kantarcioglu. Access pattern disclosure on searchable encryption: Ramification, attack and mitigation. In *Network and Distributed System Security Symposium (NDSS)*, 2012.
- [Jan] Svante Janson. Tail bounds for sums of geometric and exponential variables. <https://arxiv.org/abs/1709.08157>.
- [JJK⁺13] Stanislaw Jarecki, Charanjit Jutla, Hugo Krawczyk, Marcel C. Rosu, and Michael Steiner. Outsourced symmetric private information retrieval. In *ACM Conference on Computer and Communications Security (CCS)*, 2013.
- [KKMN09] Aleksandra Korolova, Krishnaram Kenthapadi, Nina Mishra, and Alexandros Ntoulas. Releasing search queries and clicks privately. In *Proceedings of the 18th International Conference on World Wide Web (WWW)*, pages 171–180, 2009.
- [KKNO16] Georgios Kellaris, George Kollios, Kobbi Nissim, and Adam O’Neill. Generic attacks on secure outsourced databases. In *ACM CCS*, page 1329–1340, 2016.
- [KKNO17] Georgios Kellaris, George Kollios, Kobbi Nissim, and Adam O’Neill. Accessing data while preserving privacy. *CoRR*, abs/1706.01552, 2017.
- [KNRR15] Mahmoud Abo Khamis, Hung Q. Ngo, Christopher Ré, and Atri Rudra. Joins via geometric resolutions: Worst-case and beyond. In *PODS*, pages 213–228. ACM, 2015.
- [KP16] Yael Tauman Kalai and Omer Paneth. Delegating RAM computations. In *Theory of Cryptography - 14th International Conference, TCC 2016-B, Beijing, China, October 31 - November 3, 2016, Proceedings, Part II*, volume 9986 of *Lecture Notes in Computer Science*, pages 91–118, 2016.
- [KPT20] Evgenios M. Kornaropoulos, Charalampos Papamanthou, and Roberto Tamassia. The state of the uniform: Attacks on encrypted databases beyond the uniform query distribution. In *2020 IEEE Symposium on Security and Privacy, SP 2020, San Francisco, CA, USA, May 18-21, 2020*, pages 1223–1240. IEEE, 2020.
- [KS21] Ilan Komargodski and Elaine Shi. Differentially oblivious turing machines. In *ITCS*, 2021.

- [KTH⁺19] Ios Kotsogiannis, Yuchao Tao, Xi He, Maryam Fanaeepour, Ashwin Machanavajjhala, Michael Hay, and Gerome Miklau. Privatesql: A differentially private sql query engine. *Proc. VLDB Endow.*, 12(11):1371–1384, July 2019.
- [KTM⁺19] Ios Kotsogiannis, Yuchao Tao, Ashwin Machanavajjhala, Gerome Miklau, and Michael Hay. Architecting a differentially private SQL engine. In *CIDR*, 2019.
- [LSX19] Wei-Kai Lin, Elaine Shi, and Tiancheng Xie. Can we overcome the $n \log n$ barrier for oblivious sorting? In *SODA*, 2019.
- [MG18] Sahar Mazloom and S. Dov Gordon. Secure computation with differentially private access patterns. In *CCS*, 2018.
- [MY15] Shay Moran and Amir Yehudayoff. A note on average-case sorting. *Order*, 33:23–28, 2015.
- [NKW15] Muhammad Naveed, Seny Kamara, and Charles V. Wright. Inference attacks on property-preserving encrypted databases. In *ACM CCS*, page 644–655, 2015.
- [NNRR14] Hung Q. Ngo, Dung T. Nguyen, Christopher Ré, and Atri Rudra. Beyond worst-case analysis for joins with minesweeper. In *PODS*, pages 234–245. ACM, 2014.
- [PR04] Rasmus Pagh and Flemming Friche Rodler. Cuckoo hashing. *J. Algorithms*, 51(2):122–144, May 2004.
- [PRV⁺11] Vasilis Pappas, Mariana Raykova, Binh Vo, Steven M. Bellovin, and Tal Malkin. Private search in the real world. In *Annual Computer Security Applications Conference (ACSAC)*, 2011.
- [PRZB11] Raluca Ada Popa, Catherine M. S. Redfield, Nikolai Zeldovich, and Hari Balakrishnan. Cryptdb: Protecting confidentiality with encrypted query processing. In *SOSP*, page 85–100, New York, NY, USA, 2011. Association for Computing Machinery.
- [Rou20] Tim Roughgarden. *Beyond the Worst-Case Analysis of Algorithms*. Cambridge University Press, 2020.
- [RS20] Vijaya Ramachandran and Elaine Shi. Data oblivious algorithms for multicores. <https://eprint.iacr.org/2020/947.pdf>, 2020.
- [RS21] Vijaya Ramachandran and Elaine Shi. Data oblivious algorithms for multicores. In *SPAA*, 2021.
- [SBC⁺07] Elaine Shi, John Bethencourt, T.-H. Hubert Chan, Dawn Song, and Adrian Perrig. Multi-dimensional range query over encrypted data. In *IEEE Symposium on Security and Privacy*, pages 350–364, 2007.
- [SCSL11] Elaine Shi, T.-H. Hubert Chan, Emil Stefanov, and Mingfei Li. Oblivious RAM with $O((\log N)^3)$ worst-case cost. In *ASIACRYPT*, 2011.
- [Sha86] Leonard D. Shapiro. Join processing in database systems with large main memories. *ACM Trans. Database Syst.*, 11(3):239–264, 1986.
- [SO92] Micha Sharir and Mark H. Overmars. A simple output-sensitive algorithm for hidden surface removal. *ACM Trans. Graph.*, 11(1):1–11, January 1992.

- [SPS14] Emil Stefanov, Charalampos Papamanthou, and Elaine Shi. Practical dynamic searchable symmetric encryption with small leakage. In *Network and Distributed System Security Symposium (NDSS)*, 2014.
- [Sur19] Ananda Theertha Suresh. Differentially private anonymized histograms. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *NeurIPS*, pages 7969–7979, 2019.
- [SWP00] Dawn Xiaodong Song, David Wagner, and Adrian Perrig. Practical techniques for searches on encrypted data. In *Proceedings of the 2000 IEEE Symposium on Security and Privacy*, Washington, DC, USA, 2000. IEEE Computer Society.
- [Vad17] Salil Vadhan. The complexity of differential privacy. 2017.
- [Vit01] Jeffrey Scott Vitter. External Memory Algorithms and Data Structures: Dealing with Massive Data. *ACM Comput. Surv.*, 33(2):209–271, June 2001.
- [WCM18] Sameer Wagh, Paul Cuff, and Prateek Mittal. Differentially private oblivious RAM. *PoPETs*, 2018(4):64–84, 2018.
- [WCS15] Xiao Shaun Wang, T-H. Hubert Chan, and Elaine Shi. Circuit ORAM: On Tightness of the Goldreich-Ostrovsky Lower Bound. In *CCS*, 2015.
- [Yan81] Mihalis Yannakakis. Algorithms for acyclic database schemes. In *VLDB*, page 82–94, 1981.
- [ZDB⁺17] Wenting Zheng, Ankur Dave, Jethro G. Beekman, Raluca Ada Popa, Joseph E. Gonzalez, and Ion Stoica. Opaque: An oblivious and encrypted distributed analytics platform. In *NSDI*, 2017.

Appendices

A Additional Preliminaries

A.1 Oblivious Bin Placement

We present a modified version of the oblivious bin placement algorithm by Chan and Shi [CS17] for the case when all the bins do not have the same capacity.

Recall that $m \in \mathbb{N}$ denotes the number of bins and s_1, s_2, \dots, s_m denote the capacity of the bins. In the input array \mathbf{I} , each entry is either a real element of the form (v, β) where v denotes a payload string of fixed length and $\beta \in [m]$ denotes the desired bin number; or a filler element of the form (\perp, \perp) . It is promised that at most s_β elements in \mathbf{I} want to go to the bin numbered β for $\beta \in [m]$.

1. For each bin index $\beta \in [m]$, append s_β special elements of the form (\star, β) to the resulting array. At this moment, all real elements and \star elements are tagged with β which is the element’s desired bin number. These \star elements ensure that every bin $\beta \in [m]$ will receive at least s_β elements after the next step.
2. Obviously sort the resulting array by the desired bin index, placing all fillers at the end. When elements have the same bin number, place the special \star elements after real elements.

3. In a linear scan, for each bin index $\beta \in [m]$, tag the first s_β (real or special) elements wanting to go to bin β as **normal**; and tag all remaining elements wanting to go to bin β as **excess**.
4. Obviously sort the resulting array placing all elements tagged with **excess** and all filler elements at the end; and all other elements should be ordered by their desired bin number.
5. Truncate the resulting array keeping only the first $S := \sum_{\beta \in [m]} s_\beta$ elements — the resulting array is called **O**.
6. For every element in **O**: if it is a special \star element, replace with a filler; further, remove **excess** or **normal** tags that are only needed internally by this algorithm. Output the resulting **O**.

The correctness of the above algorithm can be checked mechanically, and for obviousness, notice that the algorithm's access patterns are fully determined by $|\mathbf{I}|$ and s_1, \dots, s_m . Therefore, the algorithm satisfies obliviousness w.r.t. to the leakage $|\mathbf{I}|$ and s_1, \dots, s_m . The algorithm's runtime and cache complexity are dominated by a constant number of oblivious sortings, and the largest array sorted has length upper bounded by $|\mathbf{I}| + S$ where $S := \sum_{\beta \in [m]} s_\beta$.

A.2 Sum of Independent Geometric Random Variables

We need the following measure concentration bound for the summation of independent geometric random variables.

Theorem A.1 (Measure concentration for summation of independent geometric random variables [Jan]). *Let X be a random variable denoting the summation of n independent geometric random variables with the parameter $p \in (0, 1]$. Then, for any $\lambda \geq 1$, we have that*

$$\Pr(X \geq \lambda\mu) \leq \exp(-n \cdot (\lambda - 1 - \ln \lambda))$$

Furthermore, for any $\lambda \leq 1$, we have that

$$\Pr(X \leq \lambda\mu) \leq \exp(-n \cdot (\lambda - 1 - \ln \lambda))$$

Note that if we clamp the geometric to some upper bound U (i.e., round it down to U if it is greater than U), the above theorem still holds due to a stochastic dominance argument.

B Securely Sampling Noises on a Finite Computer

So far we have assumed sampling from the shifted and truncated geometric distribution in $O(1)$ time. We now discuss how to get rid of this assumption.

Henceforth we assume the typical parameters $\epsilon = \Theta(1)$ and $\delta = 1/N^c$. For convenience, we may assume that the sampled distribution need not be shifted so we can pretend that the center is 0 — shifting can always be performed in a post-processing step. We may also assume that $\exp(\epsilon) = 1 + 2^{-k}$ for some integer $k = O(1)$ — if this is not the case, we can round ϵ by at most a constant factor such that this expression is satisfied.

Fact B.1. *Suppose that $\exp(\epsilon) = 1 + 2^{-k}$. We can compute the binary representation of the following quantities in a streaming fashion (i.e., bit by bit) where each additional bit take $O(1)$ extra time to compute:*

- $1/\exp(\epsilon) = 1/(1 + 2^{-k})$; and

- the probability mass of the distribution $\text{Geom}(\exp(\epsilon))$ at 0, which is $2^{-k}/(2 + 2^{-k})$.

Insecure sampling algorithm. To sample from $\text{Geom}(\exp(\epsilon))$, we first sample a coin with a biased probability $p = 2^{-k}/(2 + 2^{-k})$ to decide if we should choose the outcome 0. If the outcome is not zero, we then flip an unbiased coin to decide whether we want to go left or right. Without loss of generality, assume that we want to go right, since the other direction is symmetric. Next, we keep flipping a coin of bias $\exp(-\epsilon)$ until we see heads. At this point, the number of coin flips so far is the outcome. Since we will need to truncate the geometric distribution, we need not run this algorithm for an infinite amount of time, we only need to sample at most U coins of bias $\exp(-\epsilon)$ where U is the maximum value in the support of the distribution. Sampling each biased coin needed can be accomplished using a standard trick: sample a sequence of unbiased coins, and comparing i -th sampled bits with the i -th bit in the binary representation of the bias p , until we encounter the first bit that is not equal. Given Fact B.1, the next bit in the binary representation of p can be computed on the fly in $O(1)$ time. The above algorithm takes $O(1)$ expected time to generate a biased coin.

The above naïve algorithm, however, leaks information since the time it takes to sample from the distribution is revealed.

Oblivious sampling. First, observe that for sampling a biased coin, leaking the stopping time would leak the sampled coin. However, in our case, we are sampling many coins, so we will later pad the total sampling time to some maximum value such that the probability of exceeding this time is $O(\delta)$. It is also not safe to reveal how many biased coins are sampled during the sampling of the truncated geometric. To avoid this leakage, we will sample $\Theta(\frac{1}{\epsilon} \log \frac{1}{\delta})$ number of biased coins for sampling every truncated geometric — even if the sampling process completes earlier, we will pad it to the maximum.

Using this approach, sampling from the $\mathcal{G}(\epsilon, \delta)$ distribution N times takes expected $\Theta(\frac{N}{\epsilon} \log \frac{1}{\delta})$ runtime — recall that these noises needed to mask the multiplicity of keys in either input array. Using Theorem A.1, we can show that the total runtime is upper bounded by $\Theta(\frac{N}{\epsilon} \log \frac{1}{\delta})$ except with negligible in N probability. For obliviousness, if the actual sampling completes earlier, we will still pad the runtime to $\Theta(\frac{N}{\epsilon} \log \frac{1}{\delta})$. In the final step of the algorithm, we need to sample from the $\mathcal{G}(\epsilon, \delta, \Delta)$ distribution once to mask the output length. Sampling from this distribution takes expected $\Theta(\frac{N}{\epsilon} \log \frac{1}{\delta})$ time since $\Delta = O(N + \frac{1}{\epsilon} \log \frac{1}{\delta}) = O(N)$. Using Theorem A.1, we can show that the total runtime is upper bounded by $\Theta(\frac{N}{\epsilon} \log \frac{1}{\delta})$ except with negligible in N probability. If the sampling actually completes earlier, we pad it to $\Theta(\frac{N}{\epsilon} \log \frac{1}{\delta})$ for obliviousness.

B.1 Background on Cache-Agnostic Algorithms and Cache Complexity

We provide more background on the external-memory model, cache-agnostic algorithms, and cache complexity — note that the modeling approach and metrics are standard in the algorithms literature. Interestingly, it turns out that this modeling approach a perfect match for secure processors like Intel SGX. The background text below adopts the exposition of Ramachandran and Shi [RS20] almost verbatim.

External-memory algorithms. The *two-level I/O model* in Agarwal and Vitter [AV88] is a simple abstraction of the memory hierarchy that consists of a *cache* (or *internal memory*) of size M , and an arbitrarily large *main memory* (or *external memory*) partitioned into blocks of size B . Each block can contain multiple memory words. An algorithm is said to have caused a *cache-miss* (or *page fault*) if it references a block that does not reside in the cache and must be fetched from the main memory. The *cache complexity* (or *I/O complexity*) of an algorithm is the number of block

transfers or I/O operations it causes, which is equivalent to the number of cache misses it incurs. Algorithms designed for this model often crucially depend on the knowledge of the parameters M and B , and thus do not adapt well when these parameters change — such algorithms are referred to as *cache-aware* algorithms.

Cache-agnostic algorithms. The *cache-agnostic model* (originally called the cache-oblivious model), proposed by Frigo et al. [FLPR99], is an extension of the two-level I/O model which assumes that an optimal cache replacement policy is used, and requires that the algorithm *be unaware of cache parameters M and B* . A *cache-agnostic* algorithm is flexible and portable: if a cache-agnostic algorithm achieves optimal cache complexity, it means that the number of I/Os are minimized in between any two adjacent levels in the memory hierarchy (e.g., between CPU and memory, between memory and disk) — if the algorithm is executed on a multi-level storage hierarchy. The assumption of an optimal cache replacement policy can be reasonably approximated by a standard cache replacement method such as LRU [FLPR99, Dem02].