

Range Search over Encrypted Multi-Attribute Data*

Francesca Falzon[†]
francesca_falzon@brown.edu
Brown University
University of Chicago

Evangelia Anna
Markatou[†]
markatou@brown.edu
Brown University

Zachary Espiritu
zesp@brown.edu
Brown University

Roberto Tamassia
roberto@tamassia.net
Brown University

ABSTRACT

This work addresses expressive queries over encrypted data by presenting the first systematic study of multi-attribute range search on a symmetrically encrypted database outsourced to an honest-but-curious server. Prior work includes a thorough analysis of single-attribute range search schemes (e.g. Demertzis et al. 2016) and a proposed high-level approach for multi-attribute schemes (De Capitani di Vimercati et al. 2021). We first introduce a flexible framework for building secure range search schemes over multiple attributes (dimensions) by adapting a broad class of geometric search data structures to operate on encrypted data. Our framework encompasses widely used data structures such as multi-dimensional range trees and quadrees, and has strong security properties that we formally prove. We then develop six concrete highly parallelizable range search schemes within our framework that offer a sliding scale of efficiency and security tradeoffs to suit the needs of the application. We evaluate our schemes with a formal complexity and security analysis, a prototype implementation, and an experimental evaluation on real-world datasets.

KEYWORDS

Encrypted Database; Database Reconstruction; Attack

1 INTRODUCTION

With the rise of cloud services, there is a growing need for schemes that support complex privacy-preserving queries. In this paper, we study the security of schemes that support range-queries over multi-attribute data. We consider a two party setting in which a client outsources their data to a cloud provider. The server is untrusted and assumed to be an *honest-but-curious* (or *passive persistent*) adversary, and the client wishes to query their data *privately*. This threat model is standard in the encrypted database literature and enables us to model the scenario where the adversary cannot actively interfere with the protocol, but has knowledge of past queries. One solution is to use strong cryptographic primitives like fully-homomorphic encryption [22] or oblivious RAM [24]. While they offer the best security, these solutions are not yet practical. As an alternative, solutions have been proposed using *searchable symmetric encryption* (SSE) (see, e.g., [7–11, 13, 23, 35–37, 53, 55, 61]). SSE schemes offer the following tradeoff: in exchange for efficiency they reveal some well-defined information, or *leakage*, about the queries and underlying data.

Existing efficient schemes support range queries on only single-attribute (1D) data or lack formal leakage analysis. In this paper, we present the security of a broad class of schemes that support

range queries over multi-attribute data. We first give a general framework for building such schemes and give a security proof for all schemes that fit this framework. The schemes from prior work most related to those developed within our framework are by Demertzis et al. [14, 15], who present 1D range schemes with storage and security trade-offs, and by Faber et al. [19], who build on the SSE scheme in [10] to support 1D range, substring, wild-card, and phrase queries. We extend these schemes to support multi-attribute range queries, formalize the leakage of our schemes, and evaluate our schemes using real-world datasets. Our work is the first to systematically study schemes for encrypted range search on more than two attributes.

At a high level, our framework operates by decomposing the multi-dimensional domain \mathcal{D} into subranges. Given a database over domain \mathcal{D} , each subrange can be mapped to the corresponding set of records with values in that subrange. The resulting map can then be encrypted using a standard SSE (e.g. [10]). The schemes within our framework are highly parallelizable and updates can be supported by batching updates as done in [14].

Our work systematizes the construction of private range schemes. We opted for data-independent data structures since this prevents the adversary from inferring information about the data distribution from the underlying data structure. This is a consequence of the fact that leakage mitigation and scheme efficiency are fundamentally at odds with each other; that said, our framework encompasses a wide range of classical data structures including data-dependent data structures such as *kd*-trees. We address a number of challenges that arise in multiple dimensions, i.e., the *curse of dimensionality*. Classic multi-dimension data structures are significantly more complex than their 1D variants, and extending such data structures to suit the encrypted setting requires new techniques.

Despite a large body of work on SSE-based range schemes spanning 7 years and dating back to 2015 [19], there are only a handful of papers on 2D schemes, and none in higher dimensions (see Section 1.1). We present the first concrete SSE range schemes over multiple dimensions, providing a number of efficiency and security trade-offs. We give a thorough performance and leakage analysis of our schemes, thus aiding practitioners to make informed choices in the deployment of encrypted range search based on the application context and risk assessment. Our schemes are rooted in classic data structures that are well studied, optimized, efficient, extendable to other query types, and easily implementable.

1.1 Prior and Related Work

Schemes. *Order preserving encryption* (OPE) supports range queries by enabling order comparisons of the underlying plaintexts with out decrypting. However, OPE schemes have been shown vulnerable

*This paper is to appear in VLDB 2023 with minor edits.

[†]Both authors contributed equally to this research.

to several leakage abuse attacks [4, 17, 29, 52]. SSE leaks strictly less than OPE, and a number of range schemes using SSE have also been suggested (e.g., [14, 19, 32, 66]). SSE achieves its efficiency by using light-weight cryptographic primitives like *pseudorandom functions* (PRFs) and *hash functions*.

Several SSE schemes have been developed for single-attribute (1D) range queries on encrypted databases (e.g., [32, 34, 59, 66] and see for a comparative evaluation [6]). The schemes most related to those within our framework are by Demertzis et al. [14, 15], who present 1D range schemes, and by Faber et al. [19], who extend the SSE scheme in [10] to support range, substring, wild-card, and phrase queries. Both works build a binary tree on the domain, build an index based on this tree, and then use a black-box SSE scheme to encrypt the index. Bogatov et al. [5] present a differentially private solution that combines oblivious RAM and differentially private sanitizers to mitigate leakage. Wang and Chow [65] describe a framework for range-supporting schemes, however, they do not prove different properties of different range covering schemes as we do here. They then describe two forward-secure, dynamic schemes for 1D range queries. Falzon et al. [20] sketch out what we refer to as the quadratic scheme, however, this scheme is inefficient and does not scale well to higher dimensions.

De Capitani di Vimercati et al. [16] index multi-dimensional encrypted data by recursively partitioning records into boxes. However, their approach doesn't fall into the standard SSE definition [13] and they do not provide a formal leakage analysis. Thus, we do not include [16] in our comparison with prior work. Other prior works do not use symmetric encryption, but instead use public-key cryptography [60, 64] or homomorphic encryption [40] to support multi-attribute range queries, however, these schemes leak significantly more than those built on SSE.

LEAKAGE ABUSE ATTACKS. Common types of SSE scheme leakage are access pattern (the adversary can identify the encrypted records in each response), volume pattern (the adversary can observe the number of encrypted records in each response), and search pattern (the adversary can determine if two issued queries are equal). Database reconstruction attacks have been presented against schemes supporting 1D range queries. The first such attack by Kellaris et al. [39] was followed by more efficient attacks for one-dimensional range queries using access (e.g. [27, 42, 45, 50]) and volume pattern (e.g. [26, 30, 43]). Two attack works on generic two-dimensional (2D) or two-attribute database reconstruction have also been presented [20, 49] using access and search pattern leakage.

1.2 Contributions

New challenges arise in higher dimensions due to increased structural complexity and then requires fundamentally new techniques, novel data structures that can efficiently support range queries in the encrypted setting, and leakage analysis. One such problem we address is extending the standard quadtree to a novel tree data structure with a provable reduction in the number of false positives.

We introduce a **general framework** for building private range search schemes by adapting a broad class of geometric search data structures (e.g. range trees and quadtrees) to operate on encrypted data. Our schemes reduce a range query to a set of queries on an

Table 1: Comparison of selected non-interactive range search schemes from [14] and [19] with our schemes (bold). We assume that all schemes are instantiated using an EMM or SSE scheme that hides access pattern and leaks search and volume pattern. We show the asymptotic complexity of the schemes. Notation: range size R , result size r , database size n , domain size m , number of dimensions (query attributes) d . We evaluate the security of the schemes based on token co-occurrences (reduce security) and false positives (increase security). Notation for token co-occurrences: ● high; ○ moderate; ○ none. Note that the query size for Quad-BRC is worst-case.

Scheme	Complexity				Security	
	Query Size	Resp. Size	Storage	Dim.	Token Co-occur.	False pos.
Range Tree [19]	$\log R$	r	$m + n \log m$	1	○	-
Quadratic [14]	1	r	$m + n^2$	1	○	-
Constant [14]	$\log R$	r	$m + n$	1	●	-
Log-URC [14]	$\log R$	r	$m + n \log m$	1	○	-
Log-BRC [14]	$\log R$	r	$m + n \log m$	1	○	-
Tdag-SRC [14]	1	$r + R$	$m + n \log m$	1	○	X
adqratic	1	r	$m + n^2$	d	○	-
Linear	R	r	$m + n$	d	●	-
Range-URC	$\log^d R$	r	$m + n \log^d m$	d	○	-
Range-BRC	$\log^d R$	r	$m + n \log^d m$	d	○	-
ad-BRC	$m^{\frac{d-1}{d}}$	r	$m + n \log m$	d	○	-
Tdag-SRC	1	$r + R$	$m + n \log^d m$	d	○	X
Qdag-SRC	1	$r + R^2$	$m + n \log m$	d	○	X

underlying encrypted multimap. We provide a generic security proof for *any* scheme derived from our framework (Section 3).

We then present six concrete range search schemes that fit within our framework and support queries on multiple attributes. The schemes offer a sliding scale of efficiency and security tradeoffs to suit the needs of the application (Table 1 and Section 4). One scheme reduces range queries to point queries. Three schemes are based on the range tree and build upon the 1D schemes by Demertzis et al. [14, 15] and by Faber et al. [19]. One is based based on the 2D quadtree [21], another fundamental spatial search data structure. The last scheme uses a novel data structure that is based on the quadtree and is designed to reduce the number of search tokens needed. These range-search data structures are best used for a small number of attributes, and even in the plaintext setting, alternative approaches (e.g., project to one-dimension and filter the results [3]) should be used to handle higher dimensional data.

Our schemes employ **range covering methods** that reduce the given range query to a set of precomputed queries whose answers are stored in the multimap, extending to multiple dimensions the 1D single range cover (SRC), best range cover (BRC), and uniform range cover (URC) [41] methods presented in [14, 15].

We identify a form of leakage that we call **structure pattern**, i.e., the pattern of co-occurrence of subqueries. Structure pattern leakage is inherent in schemes derived from standard multidimensional search data structures (e.g., the 1D logarithmic scheme in [14], which is derived from the range tree). We evaluate our schemes

with a theoretical complexity analysis, prototype implementation, and experimental evaluation on real-world datasets (Section 6).

We summarize our contributions as follows:

- We are the first to build a *formal general framework* for building SSE schemes that support range search on *multiple* attributes.
- We build a number of schemes that fit our framework, are based on *data-independent* data structures, and offer different *trade-offs* between bandwidth, storage and security. Our schemes are highly parallelizable and efficient in practice.
- We formally define the leakage of any scheme that fits our framework and analyze the relative security of our schemes with examples of what information an adversary can recover given scheme-specific leakage.
- We re-examine the relative security of the basic and uniform range covers, in the context of token co-occurrence.
- We support our results with a *prototype implementation*, and *experimental evaluation* on real-world databases.

Table 1 compares our concrete schemes with selected prior work. Proofs of theorems not in the main body can be found in the appendix.

2 PRELIMINARIES

Given integers a, b with $a \leq b$, let $[a] = \{1, 2, \dots, a\}$ and let $[a, b] = \{a, a + 1, \dots, b\}$. Let m_1, \dots, m_d be positive integers and $d \geq 2$. A *d-attribute database*, or a *d-dimensional database*, D is an injective mapping from a domain $\mathcal{D} = [m_1] \times \dots \times [m_d]$ to a set of records of $O(1)$ size. We denote the set of records with domain value $x = (x_1, \dots, x_d) \in \mathcal{D}$ as $D[x]$. A *d-dimensional range query* is a hyper-rectangle $[a_1, b_1] \times \dots \times [a_d, b_d]$ where $[a_i, b_i] \subseteq [1, m_i]$ denotes the range in the i -th dimension.

PRFs. A *pseudorandom function* (PRF) family is a polynomial-time computable algorithm F that takes a key $K \in \{0, 1\}^\lambda$ and an input $x \in \{0, 1\}^*$ and returns $y \in \{0, 1\}^k$ for some integer k . A PRF should be indistinguishable from random functions to any polynomial-time adversary with all but negligible probability.

2.1 Classic Range-Supporting Data Structures

In this work, we use and build upon two classic range-supporting data structures: the range tree and the quad tree.

RANGE TREE [3]. A *range tree* G_{RT} on a d -dimensional domain is a recursively defined tree. Start with a binary search tree on $[m_1]$. Each node v in this tree is associated with a binary tree on $[m_2]$. More generally, there is an edge from each vertex of the binary trees on $[m_i]$ to the root of a binary tree on $[m_{i+1}]$. A *dyadic range* is an interval with a power-of-two length $\ell = 2^k$ and its start index is $1 \bmod \ell$. A binary search tree on $[m]$ can thus be viewed as a tree whose nodes are each associated with a dyadic range in $[m]$. The source s of G_{RT} is such that $s.range = \mathcal{D}$. For a node v of a binary tree on $[m_i]$, let $v.dyadic$ denote the dyadic range in $[m_i]$ that v is associated with. Let $w.range = w_1 \times \dots \times w_d$ be the canonical range of the root w of v 's binary subtree on $[m_i]$. We have

$$v.range = w_1 \times \dots \times w_{i-1} \times v.dyadic \times [m_{i+1}] \times \dots \times [m_d] \quad (1)$$

See Figure 2 for an example of 2D range tree.

REGION QUAD TREE [21]. The quadtree recursively subdivides the domain into 2^d quadrants or orthants. Unlike the range tree, each node of the quadtree has four children, each corresponding to one of the four quadrants covering the node's range. See Figure 3 for an example of a 2D quadtree.

2.2 EMM Definition and Security Model

EMM SCHEME SYNTAX. Our range search schemes on encrypted data are built using an *encrypted multimap* (EMM) scheme in a generic manner. A *multimap* is a map that takes labels from a label space L to sets of values from a value space V i.e. $MM : L \mapsto 2^V \cup \{\perp\}$ where \perp indicates an uninitialized value. Given a multimap MM , we denote the set of values associated to label ℓ as $MM[\ell]$.

DEFINITION 1 ([11]). An *encrypted multimap scheme* (EMM) is a tuple of algorithms $\Sigma = (\text{Setup}, \text{Query}, \text{Eval}, \text{Result})$, where

- $\Sigma.\text{Setup}$: (probabilistic) takes a security parameter and a multimap, and returns a secret key and an encrypted multimap.
- $\Sigma.\text{Query}$: takes a key and label, and returns a search token.
- $\Sigma.\text{Eval}$: takes an encrypted multimap and a search token, and returns a ciphertext.
- $\Sigma.\text{Result}$: takes a key and a ciphertext, and returns a set of values.

Our label space L is the set of possible ranges over the desired domain, and the value space V is the set of possible record values, i.e. $\{0, 1\}^*$. We use the terms SSE and EMM interchangeably.

EMM SECURITY MODEL. The security of EMMs is traditionally proven using the real-ideal paradigm [11]. The definition of adaptive security for an EMM scheme Σ is parameterized by a leakage function $\mathcal{L}^\Sigma = (\mathcal{L}_S^\Sigma, \mathcal{L}_Q^\Sigma)$ which describes the exact information that a passive adversary may learn about the underlying database. In particular, \mathcal{L}_S^Σ captures the leakage at setup and \mathcal{L}_Q^Σ captures the leakage when a sequence of queries is issued. Using this security framework, we refer to an adaptively $(\mathcal{L}_S^\Sigma, \mathcal{L}_Q^\Sigma)$ -secure EMM scheme. The goal is to prove that the EMM scheme is indistinguishable from an ideal setting in which an algorithm simulates the response of the setup and query algorithms using only the leakage. Adaptive security of an EMM scheme is formally defined below.

DEFINITION 2. Let $\Sigma = (\text{Setup}, \text{Query}, \text{Eval}, \text{Result})$ be an EMM scheme and let $\mathcal{L}^\Sigma = (\mathcal{L}_S^\Sigma, \mathcal{L}_Q^\Sigma)$ be a tuple of stateful algorithms. For distinct algorithms \mathcal{A}, \mathcal{S} , and \mathcal{C} we describe two experiments below. $\text{Real}_{\mathcal{A}}^\Sigma(1^\lambda)$

- (1) The adversary \mathcal{A} selects a multimap MM and gives it to the challenger \mathcal{C} .
- (2) The challenger \mathcal{C} runs the setup algorithm with 1^λ and MM as input, $(K, EMM) \leftarrow \Sigma.\text{Setup}(MM)$. The challenger \mathcal{C} sends the encrypted multimap EMM to the adversary \mathcal{A} .
- (3) \mathcal{A} adaptively chooses labels $\ell_1, \dots, \ell_{\text{poly}(\lambda)}$; for each label ℓ_i the adversary sees the token $t_i \leftarrow \Sigma.\text{Query}(K, \ell_i)$.
- (4) \mathcal{A} eventually outputs a bit $b \in \{0, 1\}$.

$\text{Ideal}_{\mathcal{A}, \mathcal{S}}^\Sigma(1^\lambda)$

- (1) The adversary \mathcal{A} selects a multimap MM and sends MM to the challenger \mathcal{C} ; \mathcal{C} sends $\mathcal{L}_S^\Sigma(MM)$ to the simulator \mathcal{S} .
- (2) The simulator \mathcal{S} generates an encrypted multimap EMM and gives it to the adversary \mathcal{A} .

- (3) \mathcal{A} adaptively chooses labels $\ell_1, \dots, \ell_{\text{poly}(\lambda)}$; for each label ℓ_i , \mathcal{C} gives $\mathcal{L}_Q^\Sigma(\text{MM}, \ell_i)$ to \mathcal{S} , and \mathcal{S} outputs a token t_i to \mathcal{A} .
- (4) \mathcal{A} eventually outputs a bit $b \in \{0, 1\}$.

Scheme Σ is **adaptively** \mathcal{L}^Σ -secure if for all polynomial-time adversaries \mathcal{A} , there exists a poly-time simulator \mathcal{S} such that:

$$|\Pr[\text{Real}_{\mathcal{A}}^\Sigma(1^\lambda) = 1] - \Pr[\text{Ideal}_{\mathcal{A}, \mathcal{S}}^\Sigma(1^\lambda) = 1]| \leq \text{negl}(\lambda).$$

FORMALIZING LEAKAGE EMMs are parameterized by different leakage functions, which output information about the underlying data structure and its contents. Below, we define two common leakage functions of EMM schemes relevant to this work. Let MM be a multimap with label space \mathcal{L} and volume space \mathcal{V} .

- The **search pattern** reveals when two queries are equal. It takes as input a multimap MM and a label $\ell \in \mathcal{L}$, and outputs an ID. Without loss of generality we assume a 1-to-1 correspondence between range queries and identifiers: $\text{SP}(\text{MM}, \ell) \mapsto i \in [|\mathcal{L}|]$.
- The **volume pattern** of a label ℓ reveals the number of records in $\text{MM}[\ell]$: $\text{Vol}(\text{MM}, \ell) = |\text{MM}[\ell]|$.

In this paper we assume that the underlying EMM scheme is response-hiding. Attacks on multi-dimensional range queries have leveraged access and search pattern [20, 49], and we thus chose an underlying EMM scheme that does not leak access pattern.

2.3 REMM Definition and Security Model

DEFINITION 3. A **range encrypted multimap scheme** is a tuple of four algorithms $\text{REMM} = (\text{Setup}, \text{Query}, \text{Eval}, \text{Result})$. The syntax of the algorithms is defined as those in Definition 1 with the following two changes:

- REMM.Setup takes as input a security parameter 1^λ and a database D and outputs a key K and an encrypted database EMM .
- REMM.Query takes as a input a key K and a range query q on the domain of D and outputs a token t .

In order to support range queries, we take the label space of the underlying multimap to be the set of all possible range queries and the value space to be the set of all possible records.

For correctness we require that for all d -dimensional databases D with domain \mathcal{D} , all d -dimensional range queries q over \mathcal{D} , and all security parameters 1^λ , we have $\{D[x] : x \in q\} \subseteq V$, where $(K, \text{EMM}) \leftarrow \text{REMM.Setup}(1^\lambda, D)$, $t \leftarrow \text{REMM.Query}(K, q)$, $C \leftarrow \text{REMM.Eval}(\text{EMM}, t)$, and $V \leftarrow \text{REMM.Eval}(K, C)$.

REMM SECURITY MODEL. We extend the security model in [14] to encrypted multidimensional range schemes with games $\text{Real}_{\mathcal{A}}^{\text{REMM}}$ and $\text{Ideal}_{\mathcal{A}, \mathcal{S}}^{\text{REMM}}$. They are identical to the game in Definition 2 except that in step (1) the adversary selects a multi-attribute database D on domain \mathcal{D} , in step (2) the adversary selects a polynomial number of range queries on the domain \mathcal{D} , and Σ is replaced by an encrypted multi-dimensional range scheme REMM . The adaptive security of REMM schemes is defined analogously to Definition 2.

3 GENERIC FRAMEWORK

We now present a framework for building range encrypted multimap schemes from data structures for multidimensional range search based on a search DAG. This framework generalizes many

1D range schemes and captures commonly used data-structures for range search such as range-trees, kd-trees, and quad-trees. These schemes provide a variety of trade-offs with respect to query size, response size, storage, and security. We introduce the family of range-supporting data structures and explain how to build an encrypted index from a data structure in this family. This framework allows us to characterize entire classes of schemes. Furthermore, we are able to prove the properties of DAGs and range covers that are necessary to guarantee certain properties.

DEFINITION 4. A **range-supporting data structure** for a database D with domain \mathcal{D} is a pair (G, RC) , where:

- G is a connected directed acyclic graph (DAG).
- Each vertex v of G corresponds to a range on domain \mathcal{D} , which we denote as $v.\text{range}$ and refer to as a **canonical range**.
- G has a single source vertex s whose range is the entire domain, i.e., $s.\text{range} = \mathcal{D}$. For each non-sink (non-leaf) vertex v of G , we have $v.\text{range} = \bigcup_{(v,w) \in G} w.\text{range}$.
- RC , called **range covering algorithm**, is a polynomial-time algorithm that takes as input DAG G and a range query q on domain \mathcal{D} , and returns a subset W of vertices of G , called a **cover** of range q , such that the union of the canonical ranges of W includes range q , i.e., $q \subseteq \bigcup_{w \in W} w.\text{range}$.

A range-supporting data structure can be used to perform range queries by precomputing and storing the responses to all the canonical ranges of the scheme. To perform a range query q , we use the range covering function to find a cover W of q , retrieve the responses to the canonical queries for the nodes of W , and return their union as the response to q . The response to a range query q may have **false positives**, i.e., points of the database outside of range q , which must be filtered out to obtain the exact response. To avoid false positives, we use a data structure where the cover W returned by the range covering function is such that the union of the canonical ranges of W is equal to range q , i.e., $q = \bigcup_{v \in W} v.\text{range}$. The classic 1D range tree with the best range cover (Algorithm 1) is an example of a range-supporting data structure without false positives. The theorem below gives a necessary condition for a range-supporting data structure to have no false positives.

THEOREM 1. Let (G, RC) be a range-supporting data structure for a domain \mathcal{D} such that the answer to any range query has no false positives. Then, for every domain point $x \in \mathcal{D}$, there is a node v of G with canonical range $v.\text{range} = x$.

PROOF. For a contradiction, suppose there is a point $x \in \mathcal{D}$ for which such v does not exist. If $q = x$ is queried, then by Property (v) of Definition 4, RC must return some node that covers q . Algorithm RC must thus return a node $w \in V$ such that $q \subset w.\text{range}$.

For any DAG, there is a trivial range cover that returns the source node of the DAG for every query, i.e., the cover consisting of the entire domain. In general, this algorithm will cause false positives. To avoid false positives one needs to use a cover with multiple nodes. When the DAG is a tree, T , whose leaves are associated with the domain points (Theorem 1) and for each internal node v , the canonical ranges of the children of v are a partition of $v.\text{range}$, Algorithm 1, the **best range cover** (BRC), produces a cover of the query range with the minimum number of nodes (see Figure 2).

Algorithm 1: BRC(T, q, v)

```

1: // Invoked with BRC( $T, q, s$ ), where  $s$  is the root (source) of  $T$ 
2: Label  $v$  as explored
3:  $W \leftarrow \emptyset$ 
4: if  $v.range \subseteq q$  then
5:    $W \leftarrow \{v\}$ 
6: else
7:   if  $v.range \cap q < \emptyset$  then
8:     for  $(v, w) \in T$  and  $w$  is not labeled as explored do
9:        $W \leftarrow W \cup \text{BRC}(T, q, w)$ 
10: return  $W$ 

```

Algorithm 2: SRC(G, q, v)

```

1: // Invoked with SRC( $G, q, s$ ), where  $s$  is the source of  $G$ 
2: Label  $v$  as explored
3:  $cand \leftarrow null$ 
4: if  $q \subseteq v.range$  then
5:    $cand \leftarrow v$ 
6:   for  $(v, w) \in G$  and  $w$  is not labeled as explored do
7:      $\{t\} \leftarrow \text{SRC}(G, q, w)$ 
8:     if  $|t.range| < |cand.range|$  then
9:        $cand \leftarrow t$ 
10: return  $\{cand\}$ 

```

THEOREM 2. Let (T, BRC) be a range-supporting data structure whose DAG is a tree T such that

- (1) the canonical range of the leaves (sinks) of T are in 1-1 correspondence with the domain points $x \in \mathcal{D}$;
- (2) for each internal node v of T , the canonical ranges of the children of v are a non-trivial partition of the canonical range of v , i.e., $v.range = \bigcup_{(v,w) \in T} w.range$ and $\bigcap_{(v,w) \in T} w.range = \emptyset$.

Then for any range query q , the cover returned by BRC has no false positives and is unique and minimal (i.e. smallest number of nodes).

COROLLARY 1. Let (T, BRC) be a range-supporting data structure whose DAG is a tree T such that only (2) holds. For any range query q that is the union of canonical ranges of leaves, the cover returned by BRC has no false positives and is the unique minimal cover.

Observe that for queries of the same size, BRC can produce covers of different sizes (see Figure 2). Cover size may reveal some information about the location of the queried range. Kiayias et al. [41] introduce the notion of a **uniform range cover** (URC) to resolve this problem by making the size of the tokenset depend only on the size of the range, and not on its location in the domain. Let URC denote the uniform range cover algorithm for 1D ranges from [41]. It takes as input a 1D range tree G_{RT} , a range query q , and the source node of s , and returns the uniform range cover of q in G_{RT} . One way to compute URC on q is to compute the best range cover of q and then decompose the cover's nodes into their children until the desired number of nodes is reached.

When false positives are acceptable in the query answer, it may be desirable to have a cover consisting of a single node, which is accomplished by Algorithm 2, called **single range cover** (SRC).

THEOREM 3. Let (G, SRC) be a range-supporting data structure. For any range query q , the cover v returned by SRC minimizes the

GenericRS($1^\lambda, D, \mathcal{D}, \Sigma, (G, \text{RC})$)

```

Setup( $1^\lambda, D$ ) :
1: Initialize empty multimap MM.
2: for node  $w \in G$  do
3:    $\text{MM}[w.range] \leftarrow \{D[x] : x \in w.range\}$ 
4:  $(K, \text{EMM}) \leftarrow \Sigma.\text{Setup}(1^\lambda, \text{MM})$ ;
5: return  $(K, \text{EMM})$ 

Query( $K, q$ ) :
6:  $W \leftarrow \emptyset$ ;  $t \leftarrow \emptyset$ 
7:  $W \leftarrow \text{RC}(G, q)$ 
8: for  $w \in W$  do  $t \leftarrow t \cup \Sigma.\text{Query}(K, w.range)$ 
9: permute and return  $t$ 

Eval( $t, \text{EMM}$ ) :
10:  $c \leftarrow \emptyset$ ; for  $t \in t$  do  $c \leftarrow c \cup \Sigma.\text{Eval}(t, \text{EMM})$ 
11: return  $c$ 

Result( $K, c$ ) :
12:  $v \leftarrow \emptyset$ ; for  $c \in c$  do  $v \leftarrow v \cup \Sigma.\text{Result}(K, c)$ 
13: return  $v$ 

```

Figure 1: Algorithm of the generic range encrypted multimap scheme (Definition 3) for a database D with domain \mathcal{D} built from an encrypted multimap scheme Σ (Definition 1) and a range-supporting data structure (G, RC) (Definition 4).

number of domain points of the cover outside of q , i.e. the number of potential false positives.

Since many different domain-dependent data structures can be encoded as a DAG with the properties of Definition 4, we develop a generic scheme that supports range queries given any range-supporting data structure and which makes black box use of an SSE scheme. We give implementation details of the SSE scheme in Section 3.3 and we describe our generic scheme below.

Given a range-supporting data structure (G, RC) for a database D over domain \mathcal{D} , a range encrypted multimap scheme can be derived as follows. The client initializes a multimap MM and for every node v of G sets $\text{MM}[v.range] \leftarrow \{D[x] : x \in v.range\}$. The resulting multimap is encrypted using the underlying EMM scheme and outsourced to the server. To perform a range query q , the client computes $\text{RC}(G, q, s)$ to obtain a cover W of q . Using the underlying EMM scheme, the client computes search token t_w for the canonical range $w.range$ of each node $w \in W$. This set of tokens t is permuted to remove any ordering information stemming from the range covering algorithm. The set t is then sent to the server, who then retrieves the corresponding encrypted sets of records (i.e., the encrypted responses for the canonical ranges) and returns them to the client. **Generic range encrypted multimap scheme** $\text{GenericRS}(1^\lambda, D, \mathcal{D}, \Sigma, (G, \text{RC}))$ is formally described in Figure 1.

3.1 Security

Our generic range encrypted multimap scheme leaks the size of the domain and the total size of the EMM. For each query, it leaks the tokenset and the sizes of the partial responses of each token. Thus, this scheme gives rise to an additional leakage resulting from the chosen DAG and range covering scheme. We extend the notion of “partitioning” of IDs from Demertzis et al. [14], which refers specifically to the 1D range tree scheme with an underlying EMM scheme

that leaks search and access pattern, and introduce *structure pattern leakage*, a broadly applicable characterization of leakage from queries in arbitrary dimensions on a scheme instantiated with any range-supporting data structure (G, RC) and any EMM scheme.

DEFINITION 5. Let (G, RC) be a range-supporting data structure for a d -dimensional database D with domain \mathcal{D} , MM be the multimap resulting from GenericRS, and Σ be the EMM scheme. Let s be the source node of G . The *structure pattern* of a range query q is

$$\text{Str}(D, q) \mapsto (\mathcal{L}_Q^\Sigma(MM, v.range))_{v \in RC(G, q)}.$$

The leakage of GenericRS is characterized as follows.

THEOREM 4. Given an adaptively secure EMM scheme Σ that leaks search pattern and volume pattern, and a range-supporting data structure (G, RC) for a database D with domain \mathcal{D} of size m , the generic range encrypted multimap scheme GenericRS built from Σ and (G, RC) is adaptively $(\mathcal{L}_S, \mathcal{L}_Q)$ -secure, where:

$$\begin{aligned} \mathcal{L}_S(D, \mathcal{D}) &= \mathcal{L}_S^\Sigma(MM) \\ \mathcal{L}_Q(D, q^{(1)}, \dots, q^{(t)}) &= (\text{Str}(D, q^{(i)}))_{i \in [t]}. \end{aligned}$$

We prove this theorem using a standard hybrid argument. The goal is to find a sequence of games that starts in the real world and ends in the ideal world; there must be a polynomial number of such games in the sequence and they must all be indistinguishable from each other with all but negligible probability.

PROOF. We construct a stateful simulator \mathcal{S} for Setup and Query.

$\text{EMM} \leftarrow \mathcal{S}.\text{SimSetup}(1^\lambda, \mathcal{L}_S^\Sigma(MM))$

- (1) Invoke the simulator of the underlying EMM scheme on $\mathcal{L}_S^\Sigma(MM)$ to initialize an encrypted multimap EMM.
- (2) Return EMM.

$\text{Resp} \leftarrow \mathcal{S}.\text{SimQuery}(1^\lambda, (\text{Str}(D, q^{(i)}))_{i \in [t]})$

- (1) Initialize an empty set Resp .
- (2) For each $\mathcal{L}_Q^\Sigma(MM, v.range)$ such that $v \in RC(G, s, q^{(t)})$:
 - (a) Simulator \mathcal{S} uses the RC algorithm and invokes the simulator of the EMM scheme Σ on $\mathcal{L}_Q^\Sigma(MM, v.range)$, obtains the response for $v.range$ and adds it to Resp .

It remains to show that for all probabilistic poly-time adversaries \mathcal{A} , the probability $|\Pr[\text{Real}_{\mathcal{A}}^{\text{GenericRS}} = 1] - \Pr[\text{Ideal}_{\mathcal{A}, \mathcal{S}}^{\text{GenericRS}} = 1]|$ is negligibly small. We define the following two games and conclude with a hybrid argument.

Hyb0: This is identical to $\text{Real}_{\mathcal{A}}^{\text{GenericRS}}$.

Hyb1: This is identical to **Hyb0**, except that instead of invoking $\Sigma.\text{Setup}$ and $\Sigma.\text{Query}$ we invoke the simulator of the underlying EMM scheme.

$|\Pr[\text{Hyb0}] - \Pr[\text{Hyb1}]|$ is negligibly small, otherwise the security of the underlying EMM scheme would be broken with non-negligible probability. Since the distribution of **Hyb1** is identical to $\text{Ideal}_{\mathcal{A}}^{\text{GenericRS}}$ this concludes our proof.

Note that the leakage of the generic scheme is heavily dependent on the DAG and range covering algorithm used.

3.2 Choosing the Data Structure

One important choice we made is using *data-independent data structures*. This stands in contrast to the long line of work that has developed more efficient data structures for *plaintext databases*, such as R-trees (e.g. [1, 2, 31, 38, 58]) and learned indexes (e.g. [44, 48, 57]). That said, the leakage profile of the generic scheme (GenericRS) is highly dependent on the underlying DAG and range covering algorithm. Using data-dependent data structures would inherently leak more information. For example, consider a 1D database D that can be more efficiently encrypted using a range-supporting data structure whose DAG is an unbalanced binary tree that decomposes each canonical range $[i, j]$ into subranges $[i, (j - i + 1)/2]$ and $[(j - i + 1)/2 + 1, j]$. If the adversary were able to infer the tree structure, then it would additionally learn something about the data distribution of the underlying records. Many private 1D range schemes opt for balanced range trees for this reason [14, 15, 19]. We leave open the question of reconciling the development of data-dependent data structures and leakage minimization.

3.3 Implementing an EMM

We instantiate our schemes with an EMM scheme. In our implementation, we use Π_{bas} from [9]. To query a label with Π_{bas} , the client computes two per-label keys $K_1 || K_2$ which are sent to the server as the search token. For each $i \in [n]$ in increasing order, the server computes the PRF value $F(K_1, i)$ and uses it to retrieve the corresponding encrypted value. The server then uses $F(K_2, i)$ to decrypt the value which it then returns to the client. It increments i and repeats the look up until \perp is returned. We can turn this scheme into a response-hiding one with the following, simple modification: we instead use a separate key to encrypt the values. This key is kept private and the values are decrypted client-side.

Now, we can use the specific range-supporting data structure to create a plaintext multimap, which maps each canonical range to the records inside that range. This plaintext multimap is encrypted using an EMM scheme. Whenever the client issues a query, they use an appropriate range covering algorithm to determine which plaintext canonical ranges they need to query. Then, they can use the EMM scheme to compute search tokens and decrypt the result.

4 SCHEMES

We now give concrete examples of schemes that fit this framework; these schemes include generalizations of schemes presented in [14, 15] and [19]. All the schemes are instances of range encrypted multimap schemes based on range-supporting data structures presented in Section 3. The schemes support multidimensional range searches on d attributes. The first scheme we present is the Linear scheme, which provides optimal server storage at the expense of client bandwidth. In order to reduce the required bandwidth, we turn to classic data structures for geometric searching, specifically the range tree and the quadtree. Both data structures are an improvement upon the linear scheme's bandwidth, and offer trade-offs for bandwidth and storage. In order to query our data structures, we have developed the following multi-dimensional range covering techniques, extending 1D covers in novel ways: (i) Best Range Cover (BRC), which optimizes both required query bandwidth and false positives; (ii) Uniform Range Cover (URC), which in the 1D

case, Demertzis et al. [14, 15] suggest is safer than BRC at a small bandwidth expense; and (iii) Single Range Cover (SRC), which is the most secure. Finally, we develop variations of our data structures to minimize the false positive rate when a SRC is used.

For completeness, we include the multi-dimensional Quadratic scheme in Table 1. This scheme comprises of a DAG whose nodes are in one-to-one correspondence with the set of all possible range queries, together with either BRC or SRC. While this scheme offers the best possible security, it also requires quadratic storage and hence we omit it from our formal analysis.

Following the notation used throughout the paper, we denote the database with D and the domain with \mathcal{D} . We denote their sizes as $n = |D|$ and $m = |\mathcal{D}|$. The number of domain points in a query range is referred to as **range size** and denoted with R . The number database of records within a query range is referred to as **result size** and denoted with r . In our schemes, a query is issued by the client to the server as a tokenset. We refer to the number of search tokens in the tokenset as the **query size**. A response is returned by the server to the client as a collection of encrypted sets of records (one set per search token), whose total number of records is referred to as **response size**. Note that the response size is equal to the result size plus the number of false positives returned.

4.1 Linear Scheme

We first present a simple scheme, called Linear, that offers the smallest storage at the expense of the least security.

The Linear scheme indexes each record by its domain value. Intuitively, every time the client wishes to query a range q , they send a search token for every domain point of q to the server. Linear's DAG, G_L , is a star comprising a source s connected to m sinks. We have that $s.range = \mathcal{D}$ and each sink v is associated with a distinct domain point $x \in \mathcal{D}$ such that $v.range = x$. For this scheme, the generic BRC algorithm takes $O(m)$ time. To increase efficiency, we use the **linear range covering algorithm** (LRC), where in a pre-processing step, the sinks of G_L are stored in a d -dimensional array, $V[\mathcal{D}]$, indexed by its domain point. We summarize the scheme and its complexity with the following theorem.

THEOREM 5. *Let G_L be the star DAG for a database D of size n on a d -dimensional domain \mathcal{D} of size m , and let LRC be the linear covering algorithm defined above. Then (G_L, LRC) is a range-supporting data structure (Definition 4) and the range encrypted multimap scheme derived from it (Definition 5) uses space $O(n + m)$. Also, a query with range size R and result size r has query size R and response size r .*

PROOF. It is straightforward to verify properties (i) to (iv). It remains to show that algorithm LRC runs in polynomial-time when invoked on the source of G_L . LRC is implemented such that the sinks of G_L are stored in an array, $V[\mathcal{D}]$, indexed by domain point. Let q be any range query on the domain. Initializing an empty set W can be done in constant time. Looping through $x \in q$ and adding $V[x]$ to W takes time $O(R)$, where R is the size of the range q .

The linear scheme generates a multimap with m labels, one for each sink in G_L . Each record is stored once with its corresponding point value. The index has size $n + m$. When the client issues a range query of size R , the client computes R search tokens and sends them to the server. Each search corresponds to the domain points of the query and thus the response has size $O(r)$.

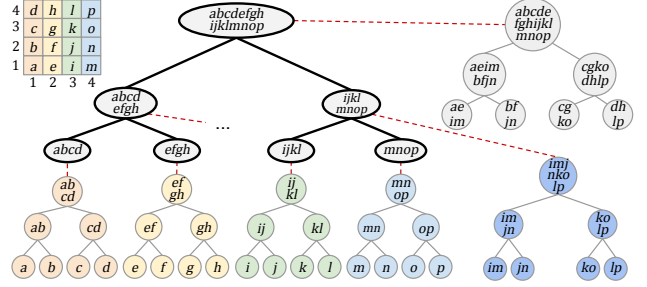


Figure 2: Example of a range tree for a database over a $[4] \times [4]$ domain. The binary tree over the first dimension is shown with black thick lines and the binary trees over the second dimension are shown with gray thin lines. Using BRC_{RT} , query range $q' = [2, 4] \times [2, 3]$ (in blue) corresponds to cover $\{v_1, v_2, w_1, w_2\}$ and query range $q'' = [3, 4] \times [2, 3]$ (in red) corresponds to cover $\{w_0, w_2\}$. The queries have the same size, but correspond to best range covers of different sizes.

4.2 Range-BRC and Range-URC Schemes

In an effort to decrease the bandwidth of the linear scheme, we present Range-BRC and Range-URC based on the range tree [3].

The multi-dimensional range tree for $d \geq 1$ does not satisfy the properties of the tree in Theorem 2. The multi-dimensional range tree can be viewed as being composed of subtrees that subdivide the domain along different dimensions; these subtrees do satisfy the properties of Theorem 2. For example, note that the three children of the root in Figure 2 correspond to the ranges $\{abcdefghijklnop\}$, $\{ijklmnop\}$, $\{abcdefghijklnop\}$ and thus do not correspond to a partition of the root's canonical range. However, the children outlined in bold correspond to the ranges $\{abcdefghijklnop\}$ and $\{ijklmnop\}$, and belong to the same subtree (which divides the domain along the first dimension). Hence they form a partition of the root's canonical range.

We thus design a **best range cover** for multi-dimensional range trees, Algorithm 3 (BRC_{RT}), that calls BRC as a subroutine on these subtrees. In particular, as we traverse the subtrees starting from the root, we compute BRC of the range along the corresponding dimension; this cover corresponds to a collection of subtrees rooted at each vertex in the cover. We then recurse and apply BRC to each of these subtrees, until we reach the final set of subtrees that partition the range along the n -th dimension. For example, in Figure 2, we have query $q'' = [3, 4] \times [2, 3]$. First, we find the best range cover in the bold tree, which corresponds to the range along the first dimension $[3, 4]$, which is node $\{ijklmnop\}$. Then, we find the best range cover in the subtree branching off of $\{ijklmnop\}$, which corresponds to the range along the second dimension $[2, 4]$, which is nodes w_0 and w_2 .

THEOREM 6. *Let G_{RT} be a multi-dimensional range tree. For any range query q , the cover returned by BRC_{RT} has no false positives and is the unique minimal cover.*

This theorem follows from standard methods in geometric search. Its proof can be found in Appendix C.

We extend the 1D URC algorithm to higher dimensions. Our uniform range cover algorithm for multi-dimensional range trees,

Algorithm 3: BRC_{RT} URC_{RT} (T, q, v)

```

1: // Invoked with input ( $T, q, s$ ), where  $s$  is the source of  $T$ 
2:  $W \leftarrow \{v\}$ 
3:  $q_1 \times \dots \times q_d \leftarrow q$ 
4: for  $i \in [d]$  do
5:    $W' \leftarrow \emptyset$ 
6:   for  $w \in W$  do
7:      $w_1 \times \dots \times w_d \leftarrow w.\text{range}$ 
8:      $\hat{q} \leftarrow w_1 \times \dots \times w_{i-1} \times q_i \times [m_{i+1}] \times \dots \times [m_d]$ 
9:     Let  $\mathcal{P} \subseteq T$  be the subtree on  $[m_i]$  rooted at  $w$ .
10:     $W' \leftarrow W' \cup \text{BRC}(\mathcal{P}, \hat{q}, w)$     $W' \leftarrow W' \cup \text{URC}(\mathcal{P}, \hat{q}, w)$ 
11:   $W \leftarrow W'$ 
12: return  $W$ 

```

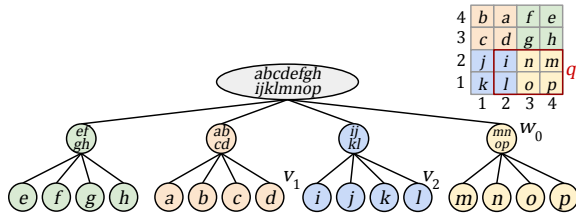


Figure 3: A quadtree for a database on domain $[1, 4] \times [1, 4]$. With BRC, query $q = [2, 4] \times [1, 2]$ corresponds to the cover $\{v_1, v_2, w_0\}$.

denoted URC_{RT} , is identical to BRC_{RT} except that on line 10, it calls URC as a subroutine instead of BRC (See Algorithm 3 for pseudocode).

THEOREM 7. Let G_{RT} be a range tree on $[m_1] \times \dots \times [m_d]$ and σ be any permutation on $[d]$. If q is a range query of size $R = R_1 \times \dots \times R_d$ and q' is a range query of size $R = R_{\sigma(1)} \times \dots \times R_{\sigma(d)}$, then their respective URC_{RT} covers W and W' are such that $|W| = |W'|$.

Recall that, given a 1D range tree G_{RT} on domain $[m]$, and range queries q and q' of the same size, then URC returns covers of the same size for q and q' [41]. Theorem 7 can be proven by extending the 1D version to multiple dimensions. We now state a theorem summarizing the complexity of the range tree scheme.

THEOREM 8. Let G_{RT} be the range tree for a database D of size n on a d -dimensional domain \mathcal{D} of size m , and let BRC_{RT} and URC_{RT} be the range covering algorithms defined in Section 4.2. We have that $(G_{RT}, \text{BRC}_{RT})$ and $(G_{RT}, \text{URC}_{RT})$ are range-supporting data structures (Definition 4) and the range encrypted multimap schemes derived from them (Definition 5) use space $O(n + m \log^d m)$. Also, a query with range size R and result size r has query size $O(\log^d R)$ and response size r .

4.3 Quad-BRC Scheme

We base our Quad-BRC scheme on the **region quadtree**. The client computes a region quadtree $G_{QT} = (V, E)$ on domain \mathcal{D} . G_{QT} is a 2^d -ary tree with one source vertex $s \in V$ such that $s.\text{range} = \mathcal{D}$. Each node $v \in V$ is associated with a d -dimensional hypercube $v.\text{range}$ and the canonical range of each child of v corresponds to one of 2^d quadrants (or *orthants*) of $v.\text{range}$. G_{QT} has m leaves, which are one-to-one with the domain values (see Figure 3). To carry out a range query q on the database the client computes the best

range cover $\text{BRC}(G_{QT}, q, s)$. We now state a theorem summarizing the complexity of Quad-BRC

THEOREM 9. Let G_{QT} be the quadtree built over d -dimensional domain \mathcal{D} of size m and let BRC be the range covering algorithm defined in Algorithm 1. We have that (G_{QT}, BRC) is a range-supporting data structure (Definition 4). Also, to store a database of size n on domain \mathcal{D} , the range encrypted multimap scheme built from (G_{QT}, BRC) according to Theorem 4 uses space $O(m + n \log m)$ and has query token size $m^{\frac{d-1}{d}}$ and ciphertext response size r , where r is the size of the plaintext query response.

4.4 Tdag-SRC Scheme

One downside of the previous schemes is that the client must often generate multiple search tokens. To overcome this, we consider *single range covers* that cover the range with only one node at the expense of false positives. Naively covering the range can result in a worst case false positive of $O(m)$ in both the range tree and quadtree constructions. Demertzis et al. propose injecting a small number of additional nodes into the range tree to create a **tree-like DAG** (TDAG) which reduces the false positive rate to $O(R)$.

To build a TDAG in one dimension, build a range tree over the domain $[m]$ and inject one extra node between every two nodes at every level of the tree. Then add edges from this node to the two nodes below it in the next level and add an edge from the node directly above in the previous level. Given a TDAG over domain $[m]$ and a range query $q \subseteq [m]$ of size R , there is a set of nodes that cover q whose canonical ranges sum to size $O(R)$ [14].

We can extend this to multiple dimensions in the following manner. Let $\mathcal{D} = [m_1] \times \dots \times [m_d]$ be the domain. Build a range tree over \mathcal{D} , then for each subtree on $[m_i]$ for $i \in [d]$, inject the nodes as described before (See Figure 4a). To issue a range query, the client uses Algorithm 2 (SRC).

COROLLARY 2. Given a TDAG constructed over the domain $\mathcal{D} = [m_1] \times \dots \times [m_d]$ and any range query $q \subseteq \mathcal{D}$ of size R , there is a set of nodes in the TDAG that cover q such that the size of their corresponding canonical ranges sum to $O(R)$.

A TDAG is only a constant factor larger than the original range trees and hence the storage requires $O(m + n \log^d m)$ space. To query this scheme, the client generates a single search token, yielding a query complexity of 1. By Corollary 2, the total number of false positives for any given query is $O(R)$ where R is the size of the query issued. We state the following theorem summarizing the complexity of the Tdag-SRC scheme.

THEOREM 10. Let G_{RS} be the TDAG built over d -dimensional domain \mathcal{D} of size m and let SRC be the single range covering algorithm defined in Algorithm 2. We have that (G_{RS}, SRC) is a range-supporting data structure (Definition 4). Also, to store a database of size n on domain \mathcal{D} , the range encrypted multimap scheme built from (G_{RS}, SRC) according to Theorem 4 uses space $O(m + n \log^d m)$ and has query token size 1 and ciphertext response size $r + R$, where r is the size of the plaintext query response and R is the size of the queried range.

4.5 Qdag-SRC Scheme

We now present Qdag-SRC, which is derived from the quadtree, supports single range covers at the expense of $O(R^d)$ false positives,

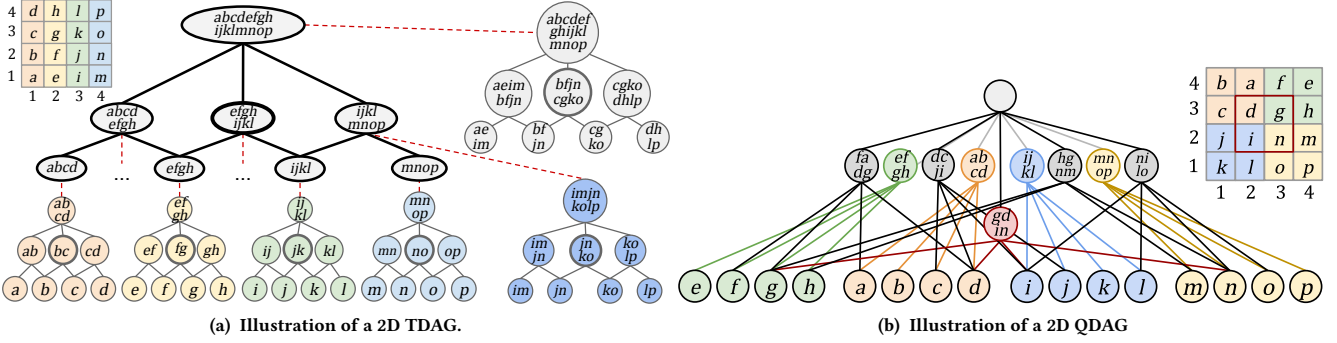


Figure 4: Examples of (a) a TDAG and (b) a QDAG for a $[4] \times [4]$ database. The inserted nodes are dark gray and the inserted edges are black.

and also extends the one-dimensional TDAG scheme [14] to higher dimensions. The response size overhead is an inherent limitation of schemes that index the domain using only hypercubes.

Given a multi-dimensional database D with domain \mathcal{D} we build the data structure the bottom up, starting with m leaves corresponding to points of domain \mathcal{D} . At level j we add nodes for all hypercubes of size 2^{n-j} tiling the domain, as well as each of these hypercubes shifted by 2^{n-j-1} along each dimension. For each node at level j , we add directed edges to all nodes in level $j-1$ which it covers. We recursively build the structure until we reach the (source) root node that corresponds to the entire domain. Each node v in this DAG is associated with a d -dimensional hypercube. To execute a range query, the client computes its SRC cover. The Qdag-SRC scheme is illustrated in Figure 4b. Below we prove that the false positive rate is $O(R^d)$ and the QDAG size is $O(m)$.

LEMMA 1. Given QDAG $G_{QS} = (V, E)$ over domain $\mathcal{D} = [m_1] \times \dots \times [m_d]$ and any range q in \mathcal{D} of size $R = R_1 \times \dots \times R_d$, there exists a vertex $v \in V$ such that $q \subseteq v.range$ and $v.range$ has size $O(R^d)$.

LEMMA 2. Let D be a database with n records and a domain \mathcal{D} of size m . Then the size of the QDAG on \mathcal{D} is $O(m)$.

We summarize below the complexity of the Qdag-SRC scheme.

THEOREM 11. Let G_{QS} be the QDAG for a database D of size n on a d -dimensional domain \mathcal{D} of size m , and let SRC be the single covering algorithm defined in Algorithm 2. We have that (G_{QS}, SRC) is a range-supporting data structure (Definition 4) and the range encrypted multimap scheme derived from it (Definition 5) uses space $O(m + n \log m)$. Also, a query with range size R and result size r has query size 1 and response size $O(r + R^2)$.

5 PERFORMANCE AND SECURITY

5.1 Performance

PARALLELIZATION. Our schemes are highly parallelizable both at setup and at query time. Our framework decomposes the domain into subranges and then constructs a map from each subrange to its corresponding records. This map can be encrypted using an EMM scheme that supports parallel setup by partitioning the label-value pairs and encrypting each batch on a different core. Search can also be parallelized in two distinct ways. Given a collection of search tokens, we can process the search for each token using a different core. Furthermore, since each search token may correspond to

multiple records, the search for each individual search token can also be parallelized, e.g. by selecting an SSE scheme that supports concurrent access requests [9, 36].

COMPLEXITY. Table 1 compares the schemes in this section. The Tdag-SRC and Qdag-SRC schemes have optimal query size but allow false positives. The other schemes avoid false positives but incur query size overhead. Notably, all the schemes have the same asymptotic storage and search time as their corresponding non-encrypted data structures. To achieve efficient client query time, the range cover algorithm builds the tokenset without instantiating the scheme's DAG, which is implicitly defined by the parameters of the domain. Thus, we can assign IDs to the nodes of the DAG so that the ID of each node in the cover is computed in $O(1)$ amortized time and space. Hence, the query time and space at the client is proportional to the query size (to generate the tokenset) plus the response size (to decrypt the received response). The query execution time at the server is also proportional to the query size plus the response size, since accessing the partial response associated with a token takes $O(1)$ expected time with an efficient multimap implementation. The client space is $O(1)$, plus $O(\log^d m)$ temporary space for the Range-BRC and Range-URC schemes, and $O(m^{\frac{d-1}{d}})$ temporary space for the Quad-BRC scheme when a query is issued, and temporary space proportional to the response size when the response is received.

EXTENSIONS. Our work specifically addresses range queries, which can be viewed as conjunctions of predicates, each associated with a 1D range query. However, our techniques support a broader class of queries involving a combination of conjunctions, disjunctions, and negations of 1D queries over different subsets of attributes. For example, the client can represent the disjunction of two 1D range queries (or the negation of a 2D range query) as the disjunction of four 2D range queries, which can be executed by taking the union of the responses to four 2D range queries. The client can use our query algorithm to find the search tokens for such queries.

Using a classic result from Boolean algebra, an arbitrary query on a d -dimensional domain involving conjunctions, disjunctions, and negations of 1D ranges can be reduced to disjunctive normal form and performed by taking the union of the responses to a collection range queries on d attributes, each computed with our methods. This collection can be identified by the client without accessing

the server. An interesting line of future work is to find an optimal query collection that minimizes the overall number of search tokens, a problem related to spatial query optimization (see, e.g., [33]). Prior work has explored SSE schemes with support for Boolean queries [10, 54]. Building such an index over each dimension to support Boolean queries over multi-attributes would result in a different leakage profile compared to that of our approach.

5.2 Security

The following theorem summarizes the leakage of our schemes.

THEOREM 12. *Let Σ be an EMM scheme that leaks search and volume pattern. The range encrypted multimap schemes Linear, Range-BRC, Range-URC, Quad-BRC, Tdag-SRC, and Qdag-SRC instantiated with Σ leak each the search, volume and structure pattern.*

We now discuss the security of our schemes. Co-occurrence stemming from access pattern has been leveraged in prior attacks [20, 49]. In Table 1, we depict the number of token co-occurrences associated with our schemes and those of prior work. We assume that all schemes are instantiated using an EMM scheme that hides access pattern and leaks search pattern and volume pattern. We consider a passive persistent (or honest-but-curious) adversary, i.e., a curious server or an eavesdropper that has compromised the communication channel. Such an adversary is strictly stronger than the **snapshot attacker** (another common academic threat model), that only gets a snapshot of the database at a particular point in time. Relevant work by Grubbs et al. [28] note that a persistent passive attacker is a better threat model for encrypted databases than a snapshot attacker, who cannot take into account past queries. Below, we extend previous approaches from attacks in 1D and 2D; these insights inform the security ranking of our schemes

LINEAR SCHEME. When a client issues a query of size R , the adversary observes R distinct search tokens, each associated with a domain point of the query. The structure pattern leaks the size of each query *and* the number of matching records. In contrast, given a scheme that leaks only access pattern, the adversary would only be able to infer the number of records matching the query. Over multiple queries, the access pattern leaks the co-occurrence of encrypted records that match each query. Structure pattern also leaks co-occurrence, since the same search token for a specific domain value can be repeatedly observed. The search tokens in this scheme are 1-1 with the domain points. Suppose we have a 2D database; if we see a response r corresponding to four search tokens, then the queried range must be of size 4. Moreover, r can either correspond to a range of size 2×2 , 4×1 , or 1×4 . If two ranges of size 4 have three overlapping search tokens, then both ranges must be of dimension 4×1 (or 1×4).

RANGE-BRC AND RANGE-URC SCHEMES. Each time the client makes a query of size R , the adversary observes $O(\log^d R)$ distinct search tokens, each token corresponding to a canonical range. The relationship between these tokens is determined by the range covering technique used. For any given range query, URC requires at least as many search tokens as BRC. For example, a range of size 2×1 under URC always results in two search tokens. Under BRC, the range may result in either one or two search tokens. This choice reveals some information about the location of the range. Specifically, if only

one search token is issued, the requested range must be a canonical range. Furthermore, canonical ranges corresponding to domain points in the inner part of the domain correspond to more queries than canonical ranges at the perimeter of the domain. A similar observation was leveraged in the attack of [20]. We conjecture that similar techniques may be applied to URC/BRC schemes.

QUAD-BRC SCHEME. The quadtree has fewer nodes than the range tree and uses less storage. A larger number of search tokens is required for each range, leaking slightly more co-occurrence information. However, the leakage is similar and we consider Quad-BRC of similar security as Range-BRC and Range-URC.

SRC SCHEMES. For the Tdag-SRC and Qdag-SRC schemes, only a single search token is sent to the server at query time. Thus they do not leak the co-occurrence of search tokens. These two schemes also allow for false positives, which make reconstruction more difficult. In fact, false positives have been introduced as a mitigation technique (e.g. [25, 51]). Note that each canonical range corresponds to a certain number of possible queries. Thus, a persistent adversary could observe a number of search tokens and exploit the frequency with which a token is observed along with knowledge of the underlying DAG structure and then try to guess the token's corresponding node. Nevertheless, we support the claim by Demertzis et al. [14] that SRC schemes are the most secure.

LEAKAGE ABUSE ATTACKS. We believe that it is possible to achieve (possibly partial) database reconstruction under the Linear scheme after observing a small number of queries, since the queries leak local information. For the tree-based BRC/URC schemes, we believe that a reconstruction attack would need more queries in order to infer the tree structure, since the leakage is less local. Finally for the SRC schemes, we believe that an attack would require a significantly larger amount of observed queries, accompanied by auxiliary information. This is due to the fact that SRC does not leak co-occurrences. It would be interesting to characterize the set of databases that produce equivalent leakage and develop attacks to better understand structure pattern.

6 EXPERIMENTS

We experimentally evaluate the performance of our schemes using the following real-world datasets (Figure 5):

GOWALLA [12]: A 4D dataset consisting of 6,442,892 latitude-longitude points of check-ins from users of the Gowalla social networking website between 2009 and 2010, a dataset used in the experiments by Demertzis et al. [14].

SPITZ [62]: A 2D dataset of 28,837 latitude-longitude points of phone location data of politician Malte Spitz from Aug 2009 to Feb 2010 and previously used in prior attack work [20, 42, 49].

CALI [47]: A 2D dataset of 21,047 latitude-longitude points of road network intersections in California, also used in [49].

NH [63]: A 3D dataset comprised of 4096 elevation points on domain $[2^6] \times [2^6] \times [2^6]$ sampled from the United States Geological Survey's Elevation Data from the White Mountains of New Hampshire. We change the domain size by keeping exactly one aggregated elevation value per latitude and longitude value.

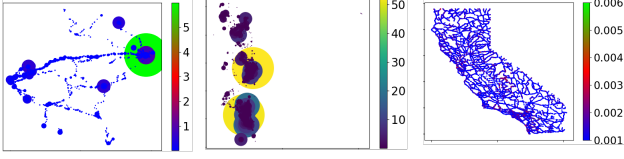


Figure 5: (Left) SPITZ, (Middle) GOWALLA, and (Right) CALI datasets on domain $2^{10} \times 2^{10}$. The color map (and relative size of the markers) indicate the number of records (in thousands) on each domain point.

IMPLEMENTATION DETAILS. We implemented our schemes in Python 3.9.2. We ran all of our experiments on a compute cluster. For simplicity, we used the same compute node for the client and the server so our results do not include any latency that would be incurred due to network transmission.

For cryptographic primitives, we used version 3.4.7 of the Python cryptography library [56]. To match the evaluation of Demertzis et al. [14], we use SHA-512 for PRFs and AES-CBC (with 128-bit block size) for encryption. We advise practitioners to follow the latest NIST standards when choosing which algorithms to use for PRFs and encryption. For our underlying EMM scheme, we used our own implementation of the Π_{bas} construction [9]. Our code is published in a public Github repository [18].

6.1 Results

In this section we use the following symbols for the schemes: Linear (\times), Range-BRC (\square), Quad-BRC (\circ), Tdag-SRC (\blacksquare), and Qdag-SRC (\bullet).

Figure 6 shows the fraction of *false positive domain points* of our SRC schemes and the number of *false positive records* they return on the CALI and SPITZ datasets. CALI has fewer false positive records returned than SPITZ since CALI is sparser, as we can see in Figure 5. Figure 7 shows the construction times for each scheme on the GOWALLA dataset from one to four dimensions. The Linear scheme remains constant since the number of nodes across the varying dimensions remains the same. In contrast, the number of canonical ranges for the SRC schemes increase as the dimension increases (as does the number of encryptions), and we thus see longer construction times. For BRC/URC, since the domain size is fixed, the trees become more shallow with increasing dimension; We conjecture that construction time decreases since less time is needed to traverse the trees.

Figures 8 and 9 show the performance of our schemes with regard to constructing, storing, and querying the index. The Linear scheme is not included in the query benchmarks due to its prohibitive complexity. Figures 8-9 give the construction, storage, and query time of each scheme on SPITZ, NH, and GOWALLA. The storage size grows with the number of database records and the domain size consistently with the asymptotic space of the schemes. The construction time also grows with the number of database records and domain size consistently with the asymptotic complexity, with the exception of the NH dataset, where large constant factors in the Qdag-SRC scheme appear to dominate the logarithmic factors of the range-based schemes for the given domain sizes. As expected, Linear has the best storage size and construction time, however, its query complexity can be prohibitive. Schemes based on the quadtree

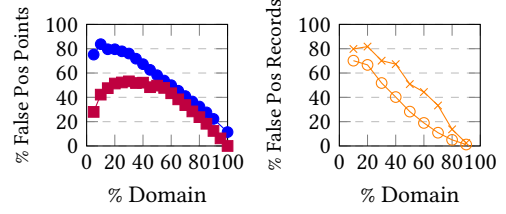


Figure 6: (Left) Percentage of *false positive domain points* by query size for the SRC schemes. The false positive rate drops off steeply for queries larger than 30% of the domain. (Right) Percent of *false positive records* returned for Qdag-SRC on the SPITZ (\times) and CALI (\circ) datasets (both on domain $2^{10} \times 2^{10}$).

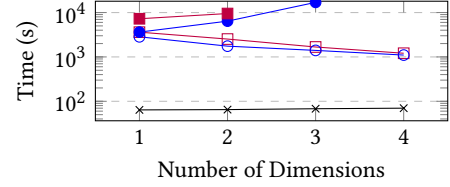


Figure 7: Construction times of our schemes on the Gowalla dataset (1 million records and 2^{24} domain points) for different dimensions.

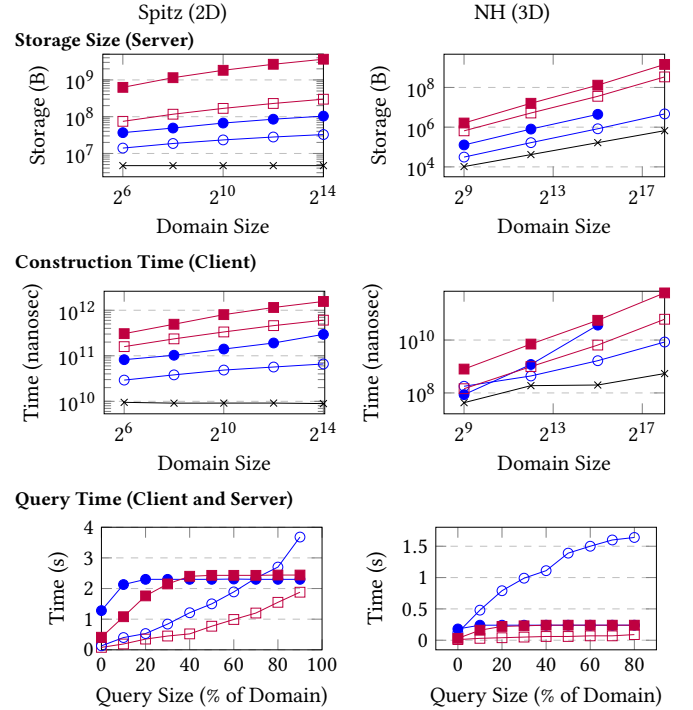


Figure 8: Storage size and construction time of the schemes by domain size on SPITZ (28,837 records) and NH (number of records $n = m^{2/3}$, where m is the domain size), with a logarithmic scale on the y -axis. Query time of the schemes by query size on SPITZ (28,837 records in domain $[2^{10}] \times [2^{10}]$) and NH (4096 records in domain $[2^6] \times [2^6] \times [2^6]$).

and QDAG have better storage size and construction time (with the exception for the NH dataset) than the schemes based on the

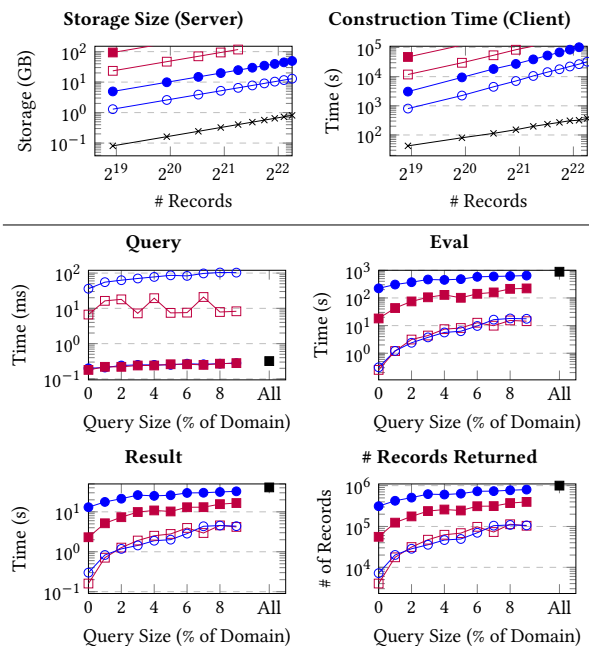


Figure 9: Storage size and construction time of the schemes by number of records on GOWALLA (domain $[2^{16}] \times [2^{16}]$). Query performance of the schemes on GOWALLA (1 million records in domain $2^{10} \times 2^{10}$) vs. a baseline (—■—) returning all records for each query. We show the times of each algorithm executed in a query (Figure 1): Query and Result at the client and Eval at the server. We also show the number of records returned. (Logarithmic scale on the y -axis)

range tree and TDAG. Quad-BRC requires the fewest nodes (with the exception of the Linear scheme) and thus its index requires less time and storage. Qdag-SRC requires a constant factor more nodes than Quad-BRC and thus requires a bit more time and storage for the index to be built. Range-BRC and Range-URC follow similarly. Finally, TDAG requires the most time and storage with a constant factor more nodes than the range tree.

The last row of Figure 8 shows the query time for SPITZ and NH, which depends on the size of range covers and the number of records that are returned. The reported query time is the sum of the time to execute methods Query, Eval, and Result (Figure 1) averaged over queries of similar size sampled uniformly at random.

For all datasets, the BRC schemes require more query time as the requested range size increases. Qdag-SRC and Tdag-SRC return the whole database for ranges larger than 20% and 40% of the domain size, respectively, and thus their query time plateaus relatively early. Note that the query time for Quad-BRC can surpasses the query time for all other schemes (e.g in the NH dataset). This is consistent with our theoretical results, as Quad-BRC’s query time complexity depends on the domain size and not on the size of the range.

The middle and last rows of Figure 9 show how the algorithms Query (generating the tokens), Eval (requesting and receiving the encrypted records) and Result (decrypting the records) perform under a database with a large number of records and small queries. We compare our approaches with a baseline of returning the whole database for every query. The Query algorithm is very efficient for both the baseline and our SRC schemes, since only one token is

returned. For the BRC schemes, Query’s runtime scales with the size of the tokenset and, as expected, Quad-BRC takes the longest, as it requires the largest tokensets. For all schemes, Query takes up less than a second, making it the most efficient function. Eval’s runtime is dominated by the number of records returned. The baseline approach takes the longest, as it needs to download the whole database. Eval in SRC schemes takes longer than in BRC schemes, and in particular, the Qdag-SRC has the greatest runtime since it has the most false positives. Result decrypts the records returned by the server, uses efficient symmetric key primitives, and scales linearly in the number of records. Our schemes return a response set that is smaller than baseline, and both BRC schemes return the same number of records. The Tdag-SRC scheme introduces some false positives and the Qdag-SRC scheme introduces the most false positives. We note that records can be processed in a streaming fashion and the client never needs to store all of them at once.

6.2 Discussion

Both in theory and practice the Linear scheme offers optimal storage complexity, however, its high query complexity makes its use prohibitive. The two SRC schemes offer the best security at the expense of higher query times (due to false positives) than the BRC schemes. Tdag-SRC and Qdag-SRC exhibit a constant factor storage overhead over the Range-BRC and Quad-BRC schemes, respectively, despite having the same asymptotic storage size.

In comparing the two BRC schemes, which have similar security, Quad-BRC uses less storage (asymptotically and experimentally) than Range-BRC. The worst-case query time for Quad-BRC, $O(m^{\frac{d-1}{d}} + r)$ is asymptotically higher than that of Range-BRC, $O(\log^d R + r)$. However, in our experiments on randomly generated queries, Quad-BRC and Range-BRC exhibit similar query performance on GOWALLA. Instead, on the SPITZ and especially the NH dataset, Quad-BRC has considerably worse query performance than Range-BRC, which is explained by the following analysis.

The ratio n/m of the number of records, n , to the domain size, m , noticeably affects query time, and this ratio varies greatly across the datasets we tested. It is 0.95 for GOWALLA, 0.0275 for SPITZ, and 0.0156 for NH. Generally, for a large ratio, query time is dominated by retrieving and decrypting the records, whereas for a small ratio, it is dominated by generating the tokenset and the number of server accesses (especially for Quad-BRC). This phenomenon is evident in Figures 8 (SPITZ and NH with small ratio) and 9 (GOWALLA with a larger ratio). For SPITZ, Quad-BRC has query time similar to the SRC schemes at query size 70% despite not returning false positives. Also, it is the slowest for larger query sizes. For NH, which has the smallest ratio, the slower query performance of Quad-BRC is even more pronounced: any queries of size greater than 10% take longer to generate the tokenset and perform the corresponding server accesses than to return the whole database with one server access. For GOWALLA, generating the tokenset takes less than a second, whereas retrieving and decrypting the records takes about 15 seconds for query size 10%. This phenomenon is not unique to the schemes we developed, but a challenge also faced by plaintext data structures. When the ratio n/m is smaller than some threshold, alternative methods should be used.

7 CONCLUSION

We introduce a framework for designing schemes that support range queries over encrypted data in multiple dimensions and describe how to turn a broad class of DAG-based spatial range search data structures into parallelizable encrypted databases that support range queries. We demonstrate the effectiveness of this framework by developing six schemes that offer trade-offs for space-complexity, query bandwidth, response size, and leakage to suit the needs of a wide variety of applications. Several aspects of our work are novel extensions of prior 1D schemes: We introduce a new scheme based on the quadtree and a new data structure that reduces the bandwidth of Quad-BRC to $O(1)$ while maintaining the same storage complexity. We adapt URC and BRC to work on the structure of the multi-dimensional range tree. The strength of our schemes lies in the fact that they are rooted in classic data structures which are efficient, flexible, and easy to implement.

ACKNOWLEDGMENTS

Work supported in part by the National Science Foundation and by the NetApp University Research Fund. Part of this research was conducted using computational resources of the Center for Computation and Visualization at Brown University.

REFERENCES

- [1] Lars Arge, Mark de Berg, Herman J. Haverkort, and Ke Yi. 2004. The Priority R-Tree: A Practically Efficient and Worst-Case Optimal R-Tree. In *Proceedings of the 2004 ACM SIGMOD International Conference on Management of Data* (Paris, France) (SIGMOD '04). Association for Computing Machinery, New York, NY, USA, 347–358.
- [2] Norbert Beckmann, Hans-Peter Kriegel, Ralf Schneider, and Bernhard Seeger. 1990. The R*-Tree: An Efficient and Robust Access Method for Points and Rectangles. In *Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data* (Atlantic City, New Jersey, USA) (SIGMOD '90). Association for Computing Machinery, New York, NY, USA, 322–331.
- [3] J. L. Bentley and J. H. Friedman. 1979. Data Structures for Range Searching. *ACM Comput. Surv.* 11, 4 (Dec. 1979), 13 pages.
- [4] V. Bindschaedler, P. Grubbs, D. Cash, T. Ristenpart, and V. Shmatikov. 2018. The Tao of Inference in Privacy-Protected Databases. *Proc. VLDB Endow.* 11, 11 (July 2018), 14 pages.
- [5] Dmytro Bogatov, Georgios Kellaris, George Kollios, Kobbi Nissim, and Adam O'Neill. 2021. epsolute: Efficiently Querying Databases While Providing Differential Privacy. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security* (Virtual Event, Republic of Korea) (CCS '21). Association for Computing Machinery, New York, NY, USA, 2262–2276.
- [6] D. Bogatov, G. Kollios, and L. Reyzin. 2019. A Comparative Evaluation of Order-Revealing Encryption Schemes and Secure Range-Query Protocols. *Proc. VLDB Endow.* 12, 8 (April 2019), 933–947.
- [7] R. Bost. 2016. Sophos: Forward Secure Searchable Encryption. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. New York, NY, USA, 12 pages.
- [8] R. Bost, B. Minaud, and O. Ohrimenko. 2017. Forward and Backward Private Searchable Encryption from Constrained Cryptographic Primitives. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security* (Dallas, Texas, USA). New York, NY, USA, 18 pages.
- [9] D. Cash, J. Jaeger, S. Jarecki, C. S. Jutla, H. Krawczyk, M.C. Rosu, and M. Steiner. 2014. 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*.
- [10] D. Cash, S. Jarecki, C. Jutla, H. Krawczyk, M.C. Rosu, and M. Steiner. 2013. Highly-Scalable Searchable Symmetric Encryption with Support for Boolean Queries. In *Advances in Cryptology – CRYPTO 2013*. Berlin, Heidelberg.
- [11] M. Chase and S. Kamara. 2010. Structured Encryption and Controlled Disclosure. In *Advances in Cryptology – ASIACRYPT 2010 – 16th International Conference on the Theory and Application of Cryptology and Information Security, Singapore, December 5-9, 2010. Proceedings (Lecture Notes in Computer Science, Vol. 6477)*.
- [12] Eunjoon Cho, Seth A. Myers, and Jure Leskovec. 2011. Friendship and Mobility: User Movement in Location-Based Social Networks. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (San Diego, California, USA) (KDD '11). New York, NY, USA, 1082–1090.
- [13] R. Curtmola, J. Garay, S. Kamara, and R. Ostrovsky. 2006. Searchable Symmetric Encryption: Improved Definitions and Efficient Constructions. In *Proc. ACM Conf. on Computer and Communications Security*.
- [14] I. Demertzis, S. Papadopoulos, O. Papapetrou, A. Deligiannakis, and M. Garofalakis. 2016. Practical private range search revisited. In *Proc. ACM Int. Conf. on Management of Data (SIGMOD)*.
- [15] Ioannis Demertzis, Stavros Papadopoulos, Odysseas Papapetrou, Antonios Deligiannakis, Minos Garofalakis, and Charalampos Papamanthou. 2018. Practical Private Range Search in Depth. *ACM Trans. Database Syst.* 43, 1, Article 2 (2018), 52 pages. <https://doi.org/10.1145/3167971>
- [16] Sabrina De Capitani di Vimercati, Dario Facchinetti, Sara Foresti, Gianluca Oldani, Stefano Paraboschi, Matthew Rossi, and Pierangela Samarati. 2021. Multi-dimensional indexes for point and range queries on outsourced encrypted data. *Proceedings of the GLOBECOM (2021)*, 1–1.
- [17] F. B. Durak, T. M. DuBuisson, and D. Cash. 2016. What Else is Revealed by Order-Revealing Encryption?. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security* (Vienna, Austria) (CCS '16). New York, NY, USA, 12 pages.
- [18] Zachary Espiritu, Evangelia Anna Markatou, Francesca Falzon, Roberto Tamassia, and William Schor. [n. d.]. *ERS*. <https://github.com/cloudsecuritygroup/ers>
- [19] S. Faber, S. Jarecki, H. Krawczyk, Q. Nguyen, M. Rosu, and M. Steiner. 2015. Rich Queries on Encrypted Data: Beyond Exact Matches. In *Computer Security – ESORICS 2015*. Cham.
- [20] F. Falzon, E. A. Markatou, Akshima, D. Cash, A. Rivkin, J. Stern, and R. Tamassia. 2020. Full Database Reconstruction in Two Dimensions. In *Proc. ACM Conf. on Computer and Communications Security (CCS)*.
- [21] R. A. Finkel and J. L. Bentley. 1974. Quad Trees a Data Structure for Retrieval on Composite Keys. 4, 1 (mar 1974), 1–9.
- [22] C. Gentry. 2009. *A Fully Homomorphic Encryption Scheme*. Ph. D. Dissertation. Stanford, CA, USA. Advisor(s) Boneh, D.
- [23] J. Ghareh Chamani, D. Papadopoulos, C. Papamanthou, and R. Jalili. 2018. New Constructions for Forward and Backward Private Symmetric Searchable Encryption (CCS '18). New York, NY, USA, 18 pages.
- [24] O. Goldreich and R. Ostrovsky. 1996. Software Protection and Simulation on Oblivious RAMs. *J. ACM* 43, 3 (May 1996), 43 pages.
- [25] Paul Grubbs, Anurag Khandelwal, Marie-Sarah Lacharité, Lloyd Brown, Lucy Li, Rachit Agarwal, and Thomas Ristenpart. 2020. Pancake: Frequency Smoothing for Encrypted Data Stores. In *USENIX Security Symposium*. 2451–2468.
- [26] P. Grubbs, M.S. Lacharité, B. Minaud, and K.G. Paterson. 2018. Pump Up the Volume: Practical Database Reconstruction from Volume Leakage on Range Queries. In *Proc. ACM Conf. on Computer and Communications Security (CCS)*.
- [27] P. Grubbs, M. Lacharité, B. Minaud, and K. G. Paterson. 2019. Learning to Reconstruct: Statistical Learning Theory and Encrypted Database Attacks. In *Proc. IEEE Symp. on Security and Privacy (S&P)*.
- [28] Paul Grubbs, Thomas Ristenpart, and Vitaly Shmatikov. 2017. Why Your Encrypted Database Is Not Secure. In *Proceedings of the 16th Workshop on Hot Topics in Operating Systems* (Whistler, BC, Canada) (HotOS '17). Association for Computing Machinery, New York, NY, USA, 162–168. <https://doi.org/10.1145/3102980.3103007>
- [29] P. Grubbs, K. Sekniqi, V. Bindschaedler, M. Naveed, and T. Ristenpart. 2017. Leakage-Abuse Attacks against Order-Revealing Encryption. In *2017 IEEE Symposium on Security and Privacy (SP)*.
- [30] Z. Gui, O. Johnson, and B. Warinschi. 2019. Encrypted Databases: New Volume Attacks against Range Queries. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS 2019, London, UK, November 11-15, 2019*.
- [31] Antonin Guttman. 1984. R-Trees: A Dynamic Index Structure for Spatial Searching. *SIGMOD Rec.* 14, 2 (jun 1984), 47–57.
- [32] F. Hahn and F. Kerschbaum. 2016. Poly-Logarithmic Range Queries on Encrypted Data with Small Leakage. In *Proceedings of the 2016 ACM on Cloud Computing Security Workshop* (Vienna, Austria) (CCSW '16). New York, NY, USA, 12 pages.
- [33] R. Helm, K. Marriott, and M. Odersky. 1995. Spatial Query Optimization: From Boolean Constraints to Range Queries. *J. Comput. System Sci.* 51, 2 (1995), 197–210. <https://doi.org/10.1006/jcss.1995.1061>
- [34] S. Kamara and T. Moataz. 2018. SQL on Structurally-Encrypted Databases. In *Advances in Cryptology – ASIACRYPT 2018*. Cham.
- [35] Seny Kamara and Tarik Moataz. 2019. Computationally Volume-Hiding Structured Encryption. In *Advances in Cryptology - EUROCRYPT 2019 - 38th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Darmstadt, Germany, May 19-23, 2019, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 11477)*. Springer, 183–213.
- [36] S. Kamara and C. Papamanthou. 2013. Parallel and Dynamic Searchable Symmetric Encryption. In *Financial Cryptography and Data Security*. Berlin, Heidelberg.
- [37] S. Kamara, C. Papamanthou, and T. Roeder. 2012. Dynamic Searchable Symmetric Encryption. In *Proceedings of the 2012 ACM Conference on Computer and Communications Security* (Raleigh, North Carolina, USA) (CCS '12). New York,

- NY, USA, 12 pages.
- [38] Ibrahim Kamel and Christos Faloutsos. 1992. Parallel R-Trees. In *Proceedings of the 1992 ACM SIGMOD International Conference on Management of Data* (San Diego, California, USA) (SIGMOD '92). Association for Computing Machinery, New York, NY, USA, 195–204.
- [39] G. Kellaris, G. Kollios, K. Nissim, and A. O'Neill. 2016. Generic Attacks on Secure Outsourced Databases. In *Proc. ACM Conf. on Computer and Communications Security*.
- [40] Shabnam Kasra Kermanshahi, Shi-Feng Sun, Joseph K. Liu, Ron Steinfield, Surya Nepal, Wang Fat Lau, and Man Au. 2020. Geometric range search on encrypted data with Forward/Backward security. *IEEE Transactions on Dependable and Secure Computing* (2020), 1–1.
- [41] A. Kiayias, S. Papadopoulos, N. Triandopoulos, and T. Zacharias. 2013. Delegatable Pseudorandom Functions and Applications (CCS '13). New York, NY, USA, 669–684.
- [42] E. M. Kornaropoulos, C. Papamanthou, and R. Tamassia. 2020. The State of the Uniform: Attacks on Encrypted Databases Beyond the Uniform Query Distribution. In *Proc. IEEE Symp. on Security and Privacy (S&P)*.
- [43] Evgenios M. Kornaropoulos, Charalampos Papamanthou, and Roberto Tamassia. 2021. Response-Hiding Encrypted Ranges: Revisiting Security via Parametrized Leakage-Abuse Attacks. In *Proc. IEEE Symp. on Security and Privacy (S&P)*.
- [44] Tim Kraska, Alex Beutel, Ed H. Chi, Jeffrey Dean, and Neoklis Polyzotis. 2018. The Case for Learned Index Structures. In *Proceedings of the 2018 International Conference on Management of Data* (Houston, TX, USA) (SIGMOD '18). Association for Computing Machinery, New York, NY, USA, 489–504.
- [45] M.S. Lacharité, B. Minaud, and K.G. Paterson. 2018. Improved reconstruction attacks on encrypted data using range query leakage. In *Proc. IEEE Symp. on Security and Privacy 2018 (S&P 2018)*.
- [46] D. T. Lee and C. K. Wong. 1977. Worst-case analysis for region and partial region searches in multidimensional binary search trees and balanced quad trees. *Acta Informatica* 9, 1 (March 1977), 23–29.
- [47] Feifei Li, Dihan Cheng, Marios Hadjieleftheriou, George Kollios, and Shang-Hua Teng. 2005. On Trip Planning Queries in Spatial Databases. In *Advances in Spatial and Temporal Databases*. Berlin, Heidelberg, 273–290.
- [48] Ryan Marcus, Andreas Kipf, Alexander van Renen, Mihail Stoian, Sanchit Misra, Alfons Kemper, Thomas Neumann, and Tim Kraska. 2020. Benchmarking Learned Indexes. *Proc. VLDB Endow.* 14, 1 (sep 2020), 1–13.
- [49] Evangelia Anna Markatou, Francesca Falzon, Roberto Tamassia, and William Schor. 2021. Reconstructing with Less: Leakage Abuse Attacks in Two Dimensions. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security* (Virtual Event, Republic of Korea) (CCS '21). Association for Computing Machinery, New York, NY, USA, 2243–2261.
- [50] Evangelia Anna Markatou and Roberto Tamassia. 2019. Full Database Reconstruction with Access and Search Pattern Leakage. In *Proc. Int. Conf. on Information Security (ISC)*.
- [51] Evangelia Anna Markatou and Roberto Tamassia. 2019. Mitigation Techniques for Attacks on 1-Dimensional Databases that Support Range Queries. In *Information Security - 22nd International Conference, ISC 2019, New York City, NY, USA, September 16-18, 2019, Proceedings (Lecture Notes in Computer Science, Vol. 11723)*. M. Naveed, S. Kamara, and C. V. Wright. 2015. Inference Attacks on Property-Preserving Encrypted Databases. In *Proc. ACM Conf. on Computer and Communications Security (CCS)*. 12 pages.
- [53] M. Naveed, M. Prabhakaran, and C. A. Gunter. 2014. Dynamic Searchable Encryption via Blind Storage. In *2014 IEEE Symposium on Security and Privacy*.
- [54] V. Pappas, F. Krell, B. Vo, V. Kolesnikov, T. Malkin, S.G. Choi, W. George, A. Keromytis, and S. Bellovin. 2014. Blind Seer: A Scalable Private DBMS. In *2014 IEEE Symposium on Security and Privacy*.
- [55] S. Patel, G. Persiano, K. Yeo, and M. Yung. 2019. Mitigating Leakage in Secure Cloud-Hosted Data Structures: Volume-Hiding for Multi-Maps via Hashing. In *Proc. ACM Conf. on Computer and Communications Security* (London, United Kingdom) (CCS '19). New York, NY, USA, 15 pages.
- [56] Python Cryptographic Authority. 2018. `pyca/cryptography`. <https://cryptography.io/> version 3.4.7.
- [57] Jianzhong Qi, Guanli Liu, Christian S. Jensen, and Lars Kulik. 2020. Effectively Learning Spatial Indices. *Proc. VLDB Endow.* 13, 12 (jul 2020), 2341–2354.
- [58] Jianzhong Qi, Yufei Tao, Yanchuan Chang, and Rui Zhang. 2018. Theoretically Optimal and Empirically Efficient R-Trees with Strong Parallelizability. *Proc. VLDB Endow.* 11, 5 (jan 2018), 621–634.
- [59] P. Rizomiliotis, E. Molla, and S. Gritzalis. 2017. REX: A Searchable Symmetric Encryption Scheme Supporting Range Queries (CCSW '17). New York, NY, USA.
- [60] E. Shi, J. Bethencourt, T-H. H. Chan, D. Song, and A. Perrig. 2007. Multi-Dimensional Range Query over Encrypted Data (SP '07). USA, 15 pages.
- [61] D. Song, D. Wagner, and A. Perrig. 2000. Practical techniques for searches on encrypted data. In *Proceeding IEEE Symposium on Security and Privacy (S&P)*.
- [62] Malte Spitz. 2011. CRAWDAD dataset spitz/cellular (v. 2011-05-04). Downloaded from <https://crawdad.org/spitz/cellular/20110504>.
- [63] United States Geological Survey. [n. d.]. Elevation Products (3DEP). Downloaded from apps.nationalmap.gov.
- [64] B. Wang, Y. Hou, M. Li, H. Wang, and H. Li. 2014. Maple: Scalable Multi-Dimensional Range Search over Encrypted Cloud Data with Tree-Based Index. In *Proc. of the 9th ACM Symposium on Information, Computer and Communications Security (ASIA CCS '14)*.
- [65] Jiafan Wang and Sherman S. M. Chow. 2022. Forward and Backward-Secure Range-Searchable Symmetric Encryption. *Proc. Priv. Enhancing Technol.* 2022, 1 (2022), 28–48. <https://doi.org/10.2478/popets-2022-0003>
- [66] Cong Zuo, Shi-Feng Sun, Joseph K Liu, Jun Shao, and Josef Pieprzyk. 2018. Dynamic searchable symmetric encryption schemes supporting range queries with forward (and backward) security. In *European Symposium on Research in Computer Security (ESORICS) (LNCS)*. Springer, 228–246.

A PROOF OF THEOREM 2

PROOF. First, we show that BRC must return an exact cover. Consider a point $x \in q$. By assumption (1) there is a leaf node $v \in V$ such that $v.range = x$. Let $p = (s, v_1, \dots, v_\ell, v)$ be the path from source s to leaf v . By assumption (2) all nodes in T that cover x must be in p .

Let u be the most recently explored vertex in p . If u satisfies the if statement on line 4, then u is an exact cover of x and is added to W . Else, we must have that $u.range \cap q < \emptyset$ and BRC is called on the children of u . Since this path starts with s then the exploration of p must start, and since p has finite length then this process must also stop. Moreover, since $v.range = x$ then the leaf v satisfies line 4 and x will necessarily be covered without false positives.

Next we show that W is minimal. Suppose not, then there is some other set $W' \subseteq V$ that is minimal. By assumption (2), if two nodes $u, v \in V$ are such that $u.range \subseteq v.range$, then u must be a descendent of v . Combining this observation, with the minimality of W' , then there must exist $v' \in W'$ and $v \in W$ such that $v.range \cap v'.range \subseteq q$. But v' must have been explored before v , and lines 4-5 would have added v' to W , which is a contradiction.

Lastly we prove uniqueness of W . Suppose there exist distinct covers W and W' of minimal size. Then there is a point $x \in q$ in the queried range that is covered by different vertices in the two covers i.e. there exist distinct $v \in W$ and $v' \in W'$ such that without loss of generality $x \subseteq v'.range \cap v.range \subseteq q$. By assumption (2), $v'.range$ forms a non-trivial partition of $v.range$. Thus there exists node v'' such that $v'.range \cup v''.range \subseteq v.range$ and W' must contain both v' and v'' instead of v , which is a contradiction.

B PROOF OF THEOREM 3

PROOF. First we demonstrate the following invariant: at the end of every iteration, $cand$ is a cover of q . We proceed inductively on the vertices explored. On the first iteration, source s is explored; s satisfies line 4 and thus $cand = s$. Since $s.range = \mathcal{D}$, then $cand$ covers q . Let v be the next vertex explored. If $q \subseteq v.range$ and $|v.range| \dot{Y} |cand.range|$, then by line 8 $cand$ is updated to v . Thus $cand$ is a cover. Otherwise $cand$ is not updated at this iteration, and by our inductive hypothesis, $cand$ must cover q .

Next we prove minimality if the number of false positives. Suppose for a contradiction, that SRC returns a cover v' of q that does not minimize the maximum number of possible false positives. Then there exists a vertex $v \in V \setminus \{v'\}$ such that $q \subseteq v.range$ and $|v.range| \dot{Y} |v'.range|$.

If v is explored before v' , then when v' is explored we have that $|cand.range| \leq |v.range|$. Since $|v.range| \dot{Y} |v'.range|$, then $cand$ would not be updated to v' , which is a contradiction.

If v is not explored before v' , then v must be explored later otherwise G would not be connected. When v is explored, $|cand.range| \leq |v'.range|$. Since $|v.range| \not\subseteq |v'.range|$, then $cand$ must be updated to a vertex whose range query is at least as small as v , which is a contradiction.

C PROOF OF THEOREM 6

PROOF. We prove the following invariant: at the end of the i^{th} iteration of the for loop on line 4, the set W is the minimal cover of the range $q_1 \times \dots \times q_i \times [m_{i+1}] \times \dots \times [m_d]$, where $q = q_1 \times \dots \times q_d$. We proceed by induction on i .

Let $i = 1$. At the start of the first iteration $W = \{s\}$. Let $\mathcal{P} \subseteq T$ be the subtree on $[m_1]$ rooted at s . Note that the canonical ranges of \mathcal{P} are comprised of dyadic ranges in dimension 1 and the whole domain in dimensions 2 to d (Equation 1). On line 8 we thus construct a query $\hat{q} = q_1 \times [m_2] \times \dots \times [m_d]$. Since the first dimension can be covered by dyadic ranges along $[m_1]$, then by Corollary 1, W contains a minimal and unique cover of \hat{q} in \mathcal{P} .

Now let $i > 1$ and w be any vertex in W . At the start of the i^{th} iteration, W must comprise of the unique, minimal cover of $q' = q_1 \times \dots \times q_{i-1} \times [m_i] \times \dots \times [m_d]$. In particular, $q' = \bigcup_{w \in W} w.range$. Let $\mathcal{P} \subseteq T$ be the subtree on $[m_i]$ rooted at w , and let $w_1 \times \dots \times w_d$ be the canonical range of w . On line 8 we construct a query $\hat{q} = w_1 \times \dots \times w_{i-1} \times q_i \times [m_{i+1}] \times \dots \times [m_d]$. Note that $\hat{q} \subseteq w.range$ and \hat{q} is the union of canonical ranges of the leaves of \mathcal{P} . By Corollary 1, Algorithm 3 (BRC_{RT}) returns the minimal and unique cover of \hat{q} in \mathcal{P} . This argument holds for all vertices $w \in W$.

If the resulting cover is not unique, then there must be another distinct cover of q' but that would either contradict the inductive hypothesis or Corollary 1. A similar argument can be made for minimality. The inductive hypothesis holds at the end of the i^{th} iteration, which concludes our proof.

D PROOF OF THEOREM 7

PROOF. We first recall that, given a 1D range tree G_{RT} on domain $[m]$, and range queries q and q' of the same size, then URC returns covers of the same size for q and q' [41].

Now let T be any of the binary subtrees of G_{RT} on $[m_i]$ for $i \in [d]$ and let w be its source. By Equation 1, the canonical ranges of the nodes in T decompose $w.range$ with respect to the dyadic ranges of $[m_i]$. More generally, for any $i \in [d]$, if q_i is a range over $[m_i]$ and T is a binary tree over $[m_i]$ where m_i is a power of 2. Then the size of a universal range cover of any range of size $|q_i|$ in T must have the same size.

Now let q be a range query of size $R = R_1 \times \dots \times R_d$ and let q' be a range query of size $R = R_{\sigma(1)} \times \dots \times R_{\sigma(d)}$ where σ is any permutation on $[d]$. For each R_i let C_i be the size of the cover URC of q_i in the binary tree over domain $[m_i]$. When we apply $URC_{RT}(T, q, s)$, the resulting cover has size $C_1 \times C_2 \times \dots \times C_d$. Similarly, cover $URC_{RT}(T, q', s)$ has size $C_{\sigma(1)} \times C_{\sigma(2)} \times \dots \times C_{\sigma(d)}$.

Thus q and q' have the same size URC covers.

E PROOF OF THEOREM 8

PROOF. We give the proof for the range encrypted multimap scheme built from (G_{RT}, BRC_{RT}) . It is straightforward to check that

DAG G_{RT} satisfies properties (i) to (iv) of Definition 4. We now show that BRC_{RT} runs in poly-time when called on G_{RT} .

At the start of the i^{th} iteration of the for loop on line 4, W contains at most $\prod_{j=1}^{i-1} \log R_j$ nodes where $R = R_1 \times \dots \times R_d$ is the size of the queried range. Parsing $w.range$ and computing \hat{q} can be done in constant time. Computing the subtree rooted at w takes time linear in the number of nodes of $\mathcal{P} \subseteq G_{RT}$. Next, observe that BRC does a depth-first search traversal of the vertices of the input tree, thus this subroutine takes time linear in the number of nodes of \mathcal{P} . When $m_i = O(m)$, the outer for loop takes time $O(m \log^d m)$ and therefore BRC_{RT} runs in polynomial time.

We now prove the complexity. The range tree has $\prod_{i=1}^d 2m_i$ nodes that each correspond to a label in the multimap; each record is stored $\prod_{i=1}^d \log m_i$ times for a total storage of $(m + n \log^d m)$. Algorithm BRC guarantees that any range can be covered by a logarithmic number of pre-computed ranges [14] in a range tree; Since we apply BRC once for every every dimension, any range over \mathcal{D} can be covered with $O(\log^d R)$ canonical ranges. Since BRC_{RT} iteratively applies BRC on one-dimensional range trees, then applying Theorem 2, the cover is minimal and has no false positives. Moreover, the cover returned by BRC on a one dimensional range tree is disjoint, and hence the cover returned by BRC_{RT} is also disjoint. Thus the response size is $O(r)$.

The proof for the range encrypted multimap scheme built from (G_{RT}, URC_{RT}) follows a similar argument and is thus omitted.

F PROOF OF THEOREM 9

PROOF. The region quadtree G_{QT} satisfies properties (i) to (iv) of Definition 4. We now show that BRC runs in polynomial time when invoked on G_{QT} and its source s . BRC implements a depth first search on G_{QT} , and stops exploring a particular path when one of two conditions is satisfied: the node v being explored is such that either $v.range \subseteq q$ (i.e. v is in the cover), or $v.range \not\subseteq q$ (i.e. v cannot be in the cover and, by property (iv) of Definition 4, neither can its children).

Else, if there is a non-empty intersection $v.range \cap q < \emptyset$ without containment, then v is not in the cover, but its descendants are and BRC is recursively called on its children. In the worst case, the client completely traverses G_{QT} which takes time $O(m)$.

We now prove the complexity of the resulting database. Given a d -dimensional domain \mathcal{D} of size m , the corresponding quadtree has $(2^d m - 1)/(2^d - 1)$ nodes and height $(\log m)/d$. Each record is stored $(\log m)/d$ times. In particular, each record is stored at the nodes on the path from the leaf associated with the record to the root. Each node corresponds to a label in the multimap and so the total storage is $O(m + n \log m) = (m 2^d - 1)/(2^d - 1) + (n \log m)/d$ where n is the total number of records.

In the worst case, the client may need to generate $O(m^{(d-1)/d})$ search tokens [46]. Since G_{QT} satisfies the properties described in Theorem 2, the cover returned by BRC has no false positives. Moreover, for any node v in G_{QT} , the canonical ranges of its children must partition $v.range$. Thus the cover is disjoint and hence each matching record is returned once for a response size of $O(r)$.

G PROOF OF THEOREM 10

PROOF. The TDAG G_{RS} satisfies properties (i) to (iv) of Definition 4. We now show that SRC runs in polynomial time when invoked on G_{RS} and its source s . SRC explores the graph in a depth-first search manner and each node is visited at most once. If $q \subseteq v.range$, then $cand$ is updated to v . Then for every vertex w such that $(v, w) \in E$, SRC is recursively called on G_{RS} and w , and a new vertex t is obtained. If $|t.range| \dot{Y} |cand.range|$, then $cand$ is updated to t . G_{RS} has $O(\log^d m)$ nodes, so visiting every node takes time $O(\log^d m)$; updating $cand$ takes constant time.

We now prove the complexity of the resulting database. Given a d -dimensional domain \mathcal{D} of size m , the corresponding range tree has $\log^d m$ nodes. Since we insert a constant number of additional nodes to each one-dimensional subtree, the overall complexity of the TDAG is also $\log^d m$. We thus have overall storage complexity

$O(m + n \log^d m)$ Since SRC returns one node as a cover, the query size is $O(1)$. By Corollary 2, the response size is $O(R + r)$.

H PROOF OF THEOREM 11

PROOF. Qdag-SRC satisfies properties (i) to (iv) of Definition 4. The proof that the range cover algorithm SRC runs in poly-time when called on G_{QS} and its source node s is analagous to the proof for SRC on TDAGs, and is omitted.

We now prove the complexity. By Lemma 2, the number of nodes in G_{QS} is a constant factor larger than the corresponding quadtree and thus the corresponding range encrypted multimap scheme has an asymptotic storage complexity of $O(m + n \log m)$. To query this scheme, the client generates a single search token, and thus the scheme has a query complexity of 1. By Lemma 1, the total number of false positives for any given query is $O(R^d)$, where R is the size of the query issued. Thus, the total response size is $O(r + R^2)$.