

Hu-Fu: Efficient and Secure Spatial Queries over Data Federation

Yongxin Tong¹ · Yuxiang Zeng¹ · Yang Song¹ · Xuchen Pan¹ ·
Zeheng Fan¹ · Chunbo Xue¹ · Zimu Zhou² · Xiaofei Zhang³ · Lei
Chen⁴ · Yi Xu¹ · Ke Xu¹ · Weifeng Lv¹

Received: date / Accepted: date

Abstract Data isolation has become an obstacle to scale up query processing over big data, since sharing raw data among data owners is often prohibitive due to security concerns. A promising solution is to perform secure queries over a federation of multiple data owners leveraging secure multi-party computation (SMC) techniques, as evidenced by recent federation studies on relational data. However, existing solutions are highly inefficient on spatial queries due to excessive secure distance operations for query processing and their usage of general-purpose SMC libraries for secure operation implementation. In this paper, we propose Hu-Fu, the first system for efficient and secure spatial query processing on a data federation. Hu-Fu seamlessly supports five mainstream spatial queries at scale, while ensuring both data and query privacy (*i.e.*, sensitive spatial information of data owners and query users). The idea is to decompose the secure processing of a spatial query into as many plaintext operations and as few secure opera-

tions as possible, where fewer secure operators are involved and all of them are implemented dedicatedly. As a working system, Hu-Fu supports not only query input in native SQL, but also heterogeneous spatial databases (*e.g.*, PostGIS, GeoMesa, and SpatialHadoop) at the backend. Extensive experiments show that Hu-Fu usually outperforms the state-of-the-arts in running time and communication cost while guaranteeing security.

Keywords Federated database · Spatial database · Query processing · Data privacy

1 Introduction

Efficient processing of spatial queries over large-scale data is essential for a wide spectrum of smart city applications, such as taxi-calling [66] and logistics planning [68]. Although the volume of spatial data continues to grow, it becomes increasingly difficult for these applications to take full advantage of the big spatial data due to the data isolation problem (*a.k.a.* isolated data) [43, 50, 52]. Spatial datasets at city or nation scale are often privately possessed and separately owned by multiple parties, where sharing raw data among parties or uploading raw data to a third party (*e.g.*, a cloud) is prohibitive due to legal regulations (*e.g.*, GDPR [55]) or commercial reasons.

A promising paradigm to tackle the data isolation problem is to perform *secure* queries over a *data federation* [14], which consists of multiple data owners (*a.k.a.* data silos) who agree on the same schema and manage their own data autonomously. Note that this paradigm differs from conventional federated databases [45] in the extra security requirement. In general, secure query processing over data federation can be solved by well-known techniques such as secure multi-party

Y. Tong E-mail: yxtong@buaa.edu.cn ·
Y. Zeng E-mail: yxzeng@buaa.edu.cn ·
Y. Song E-mail: songyangbuaa@buaa.edu.cn ·
X. Pan E-mail: panxuchen@buaa.edu.cn ·
Z. Fan E-mail: fanzh@buaa.edu.cn ·
C. Xue E-mail: xuechunbo@buaa.edu.cn ·
Z. Zhou E-mail: zimuzhou@cityu.edu.cn ·
X. Zhang E-mail: xiaofei.zhang@memphis.edu ·
L. Chen E-mail: leichen@cse.ust.hk ·
Y. Xu E-mail: xuy@buaa.edu.cn ·
K. Xu E-mail: kexu@buaa.edu.cn ·
W. Lv E-mail: lwf@buaa.edu.cn

¹ State Key Laboratory of Complex & Critical Software Environment and Beijing Advanced Innovation Center for Future Blockchain and Privacy Computing, Beihang University

² School of Data Science, City University of Hong Kong

³ The University of Memphis

⁴ The Hong Kong University of Science and Technology (GZ) and The Hong Kong University of Science and Technology

computation (SMC) [24]. Yet, only recently did pioneer studies such as SMCQL [14] and Conclave [56] take the first step towards practice with efficient query execution plans upon SMC libraries for (relational) data federation. Unsurprisingly, more applications are being built on federations of spatial data owners.

Example 1 AMAP [3] (GaoDe Map in China) has united over 8 Chinese travel companies into an integrated taxi-calling platform to offer users the taxis resources from all participating companies. A spatial data federation can protect the distribution of taxis' locations of each company (*i.e.*, data silo), which could be a business secret, from leaking to others. This privacy concern for data silos is commonly referred to as *data privacy* [28].

Example 2 During COVID-19, several mobile network operators (*e.g.*, China Mobile [4] and China Telecom [5]) cooperated as a data federation to identify individuals who had contacts with infectious patients through their location data [6]. Executing spatial queries (*e.g.*, range query) over a data federation helps identify contacts of infectious patients across multiple organizations' spatial data without compromising privacy. Here, strict privacy requirements go beyond the data privacy since the locations of patients, which appear in queries, also need protections. The privacy concern for spatial data in queries is referred to as *query privacy* [28].

Due to legal regulations (*e.g.*, GDPR [55]), protecting *data privacy* is now common in real-life scenarios, especially when spatial location implies the travel patterns or personal trajectories of a user. *Query privacy* is equally important, but perhaps gets less attention in existing research on data federation. Query privacy also has numerous real-world applications [17, 29], such as navigation, location-based social networking, location-based advertising, and POI search ¹.

Nevertheless, directly adapting the state-of-the-art data federation solutions [14, 56] to spatial data can be inefficient. From our empirical study (Sec. 2.2) of a kNN query on a real dataset, they are at least $142\times$ slower, and have at least $1,216\times$ higher communication cost than plaintext query processing. There are two reasons for such inefficiency. (i) Existing solutions process spatial queries with excessive secure distance operations, which occupy over 90% of the time cost. For example, SMCQL [14] and Conclave [56] would securely sort spatial objects by distances to the query point and pick the top-k objects, where each sorting involves numerous secure distance comparisons. (ii) Previous studies [14, 56] are built on general-purpose SMC libraries, which may

sacrifice the efficiency of specific operations for other considerations. For example, our experiment shows that the secure summation in OblivM [42], the SMC library adopted by SMCQL [14], can be accelerated by $15\times$ via dedicated implementations [23].

In this paper, we aim at efficient and secure spatial queries over a data federation, which we call *federated spatial queries*. We mainly study five queries (federated range query, range counting, kNN query, distance join, and kNN join) commonly seen in spatial database research [21, 64] and follow the semi-honest adversary model adopted by previous work [14, 56, 59]. Moreover, we develop a more practical solution than [14, 56] by eliminating the need for an honest broker and supporting more data silos (these studies support at most three data silos whereas we tested up to ten).

To this end, we propose Hu-Fu [7], a system for efficient and secure processing of federated spatial queries. As explained above, secure operations are usually slow and easily become the efficiency bottleneck. Thus, the **key idea** of Hu-Fu is to decompose a federated spatial query into as many plaintext operations while minimizing secure distance-related operations without compromising privacy. The decomposition aims to achieve two goals: (i) reduce the number of distance-related operations to the minimum, and (ii) implement secure operations faster than those in general-purpose SMC libraries. To realize this idea and implement a practical system, Hu-Fu consists of three components: an query rewriter with novel decomposition plans, a set of drivers adaptable to heterogeneous databases and an easy-to-use query interface with SQL support. Specifically, the query rewriter identifies a set of plaintext and secure operators for the query execution plan to handle the queries of interest. It ensures diverse privacy requirements, as explained in Examples 1 and 2: data privacy only, or both data and query privacy. The drivers provide implementations of secure operators with dedicated SMC protocols and plaintext operators as interfaces on top of the heterogeneous spatial databases adopted by different data silos. The query interface supports spatial queries in native SQL for easy usage.

Contribution. Our main contributions and results are summarized as follows.

- To the best of our knowledge, Hu-Fu is the first system on efficient and secure spatial queries over a data federation, and is also available on GitHub [7].
- We devise novel decomposition plans for federated spatial queries. After decomposition, an execution plan involves only a limited number of secure operators that can be effectively supported with fast and dedicated implementations.

¹ Please refer to Appendix A for more detailed explanations on these application scenarios

- Hu-Fu is an efficient, easy-to-use system that supports query input in SQL and heterogeneous spatial databases, *e.g.*, PostGIS [10], Simba [64], GeoMesa [27], SpatiaLite [11], and SpatialHadoop [21].
- Extensive evaluations show that Hu-Fu usually outperforms the state-of-the-arts [14, 56] in efficiency. Compared with two strong baselines, namely SMCQL-GIS and Conclave-GIS, which are extended from SMCQL [14] and Conclave [56] to spatial queries, Hu-Fu is up to 4 orders of magnitude faster and 5 orders of magnitude lower in communication overhead than SMCQL-GIS and Conclave-GIS with the same security level.

Compared with the preliminary version [51] of this work, we have made the following new contributions. *(i)* We expand our scope to a new and challenging setting where both data and query privacy must be preserved. Hu-Fu also provides the corresponding SQL *query interface*. *(ii)* The *query rewriter* is extended and optimized to handle all five spatial queries in this new setting. *(iii)* In *drivers*, two additional secure operators are tailored to fulfill the extra privacy requirement. *(iv)* Extensive evaluations are conducted to show the performance.

Roadmap. In the rest of this paper, we define our problem scope and identify the inefficiency of existing solutions in Sec. 2. We present an overview of Hu-Fu in Sec. 3 and elaborate on the three functional components in Sec. 4, Sec. 5, and Sec. 6. Finally, we present the evaluations in Sec. 7, review the related work in Sec. 8, and conclude in Sec. 9.

2 Problem Statement

This section clarifies our problem scope and highlights the technical challenges when developing Hu-Fu.

2.1 Problem Scope

A data federation F (“federation” as short) consists of n data silos $\{F_i\}$ (“silos” as short), where each silo holds massive *spatial objects*. Each spatial object o has a location l_o and (optionally) other attributes. The federation supports *federated spatial queries* over the spatial objects of all silos under the following settings.

- **Spatial Queries.** The federation supports mainstream spatial queries like range query, range counting, kNN query, distance join, and kNN join [44, 64].
 - **Autonomous Databases.** Each data silo is an autonomous database that manages (*e.g.*, deletes and inserts its own spatial objects and prohibits sharing its spatial objects in plaintext with the other data silos [14–16, 56].
 - **Semi-honest Adversaries.** Each silo honestly executes queries received and returns authentic results, but may attempt to infer data from other silos during query execution. This assumption is common in query processing over a data federation [14–16, 56].
- Moreover, the query processing methods should consider the following requirements thoughtfully.
- **Efficiency Requirements.** We care about the *running time* and *communication cost* to execute *exact* queries over multiple silos. Short running time is often desirable since real-life applications may process massive queries and expect prompt responses. Minimal communication cost is critical in distributed query processing [45] and secure query processing [28]. Approximate query processing over data federation [16, 20, 69] is out of our scope because applications such as contact tracing require accurate results. We consider multiple silos as aligned with real-world applications. Similar to existing solutions [14, 56], the storage efficiency, which mainly depends on silos themselves, is not our primary concern.
 - **Privacy Requirements.** We consider *two* different types of privacy requirements [24, 28].
 - (1) *Data privacy*: each data silo should not deduce any sensitive data from others, and no additional sensitive data should be revealed to the query user, except for the final query answer.
 - (2) *Query privacy* (optional): the spatial location of a user’s query cannot be revealed to data silos.

Remark. In practice, the need for query privacy may vary across applications. For example, in scenarios such as a passenger requesting a taxi-calling service through AMAP [3], query privacy may be unnecessary. This is because the platform ultimately needs to know the passenger’s pickup location. Conversely, in situations like performing contact tracings based on a patient’s location, query privacy becomes crucial to prevent the disclosure of sensitive spatial information to data silos. For ease of presentation, we classify federated spatial queries into two kinds: *asymmetric queries* (which require data privacy only) and *symmetric queries* (which require both data privacy and query privacy).

To satisfy the privacy requirements, some existing systems [14, 56] assume the existence of a trusted broker responsible for collecting partial answers, which may contain sensitive data, from each silo. In reality, even if brokers (*e.g.*, Acxiom [2]) charge high fees for their data broker services, there is still a risk of them leaking sensitive data for personal gain [1]. Thus, we explore solutions without reliance on a trusted broker.

2.2 Main Challenges

Federated queries can be realized by secure multi-party computation (SMC) [24], as in prior studies for relational data [14, 56]. Nevertheless, our empirical study shows that they are highly inefficient on spatial queries.

2.2.1 Inefficiency on Federated Spatial Queries

As an illustrative study, we perform an asymmetric federated kNN query by extending SMCQL [14] and Conclave [56], two representative solutions to secure query processing on (relational) data federations.

Overview of Existing Solutions. The common framework [14, 56] for secure query processing over a data federation decomposes query execution into *plaintext* queries within each silo and *secure* computations of the partial results across silos. Existing solutions differ in the SMC techniques used for secure operations, with garbled circuits (GC) and secret sharing (SS) as the mainstreams [24]. For example, SMCQL-GIS [14] uses a prevalent GC based library, OblivM [42], to support two silos. Conclave-GIS [56] adopts an SS based technique (Sharemind [18]), which enables query processing on three silos.

Setup. SMCQL-GIS [14] and Conclave-GIS [56] are extended to asymmetric federated kNN queries as follows. Following the “plaintext + secure” processing pipeline, each silo first conducts a plaintext kNN query and returns the k nearest points (along with their distances) to the query point. Then, the final kNNs are derived by a top- k operation from these returned points, which are securely sorted by their distances to the query point. We experiment with two silos with $k = 16$. Other details of experimental setups are elaborated in Sec. 7.1.

Result. Fig. 1 plots the (average) running time and communication cost to process an asymmetric federated kNN query leveraging existing solutions [14, 56]. The results are averaged over 50 queries. Compared with Public, *i.e.*, plaintext kNN query execution without any privacy protection, the secure counterpart incurs $142\times$ to $212\times$ longer running time and $1,216\times$ to $22,510\times$ higher communication cost. Although the method SMCQL-GIS yields shorter running time and lower communication overhead than Conclave-GIS, it is *limited to scenarios with only two silos* due to its reliance on garbled circuits (GC). Yet it still takes 2.86 seconds to answer a federated spatial query, which can hurt user experiences in applications where query time efficiency is critical.

Table 1: Percentage of time spent for *plaintext* or *secure* operations in an asymmetric federated kNN query.

Existing Solution	Plaintext	Secure
SMCQL-GIS [14]	0.14%	99.86%
Conclave-GIS [56]	0.10%	99.90%

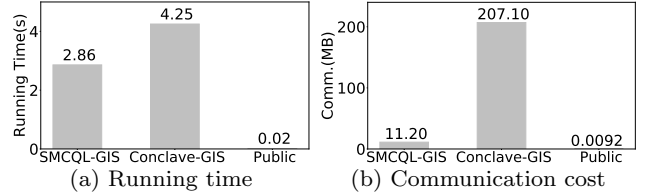


Fig. 1: Inefficiency of Conclave-GIS and SMCQL-GIS on asymmetric federated kNN query, where SMCQL-GIS and Conclave-GIS are our extensions on SMCQL [14] and Conclave [56] to spatial queries (see Sec. 7.1).

2.2.2 Understanding the Efficiency Bottleneck

Prior studies are inefficient on federated spatial queries for the following reasons.

- **Excessive Secure Distance Operations.** When processing the test query, over 99% time is spent on secure operations (*e.g.*, secure distance comparisons) as shown in Table 1. Specifically, SMCQL-GIS and Conclave-GIS adopt sorting to find kNNs among nk candidates by using $O(nk \log(nk))$ secure distance comparisons. A single secure distance comparison in SMCQL-GIS takes 209 ms, while in Conclave-GIS it takes 248 ms, which equals the time required for at least 10^6 plaintext comparisons.
- **Reliance on General-Purpose Libraries.** Existing methods use general-purpose libraries to implement secure operations (*e.g.*, OblivM [42] in SMCQL [14]). General-purpose libraries sometimes sacrifice efficiency for generalization or compatibility. For example, the secure summation used in Hu-Fu can be $16\times$ faster than that in OblivM (see Sec. 7). As will be shown in Sec. 4, federated spatial queries can be processed with only a few secure operations. This facilitates acceleration by dedicated protocols specifically tailored for these secure operations.

Takeaway. Our study shows that existing secure query processing solutions (*e.g.*, [14, 56]) for data federations are inefficient for spatial queries. The inefficiency comes from (i) massive secure distance operations, and is exacerbated by (ii) adopting general-purpose libraries for these SMC operations. In response, we propose Hu-Fu, a solution with (i) a novel execution plan for federated spatial queries that involve notably fewer secure oper-

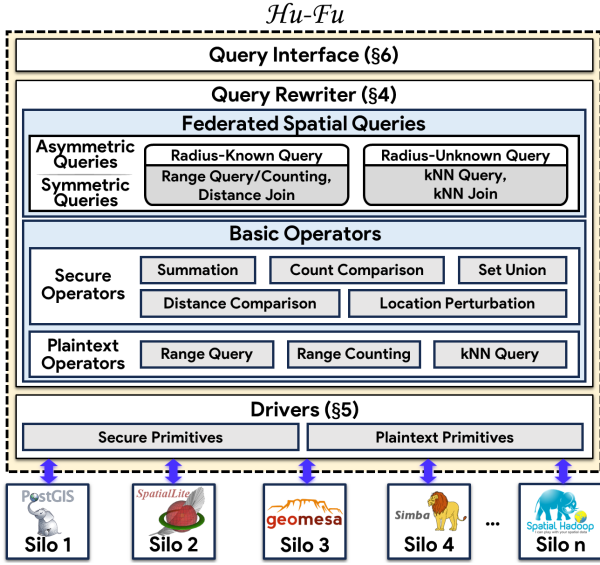


Fig. 2: Foundation architecture of Hu-Fu.

ations (see Sec. 4) and (ii) each secure operator can be implemented in high efficiency via dedicated methods (see Sec. 5). As next, we give an overview of Hu-Fu and elaborate on its functional modules in the following.

3 Hu-Fu Overview

Hu-Fu is a solution that enables efficient and secure spatial queries over a data federation. It addresses the inefficiency of federated spatial query processing via two modules: (1) a novel *query rewriter* that decomposes federated spatial queries into *plaintext* and *secure operators*, with the former executed within each silo and the latter across silos; (2) *drivers* that implement these operators as *plaintext* and *secure primitives* leveraging dedicated algorithms and optimizations. Hu-Fu also contains a transparent *query interface* to support federated spatial queries written in native SQL. We will briefly explain its architecture and workflow as follows.

3.1 Architecture

Fig. 2 shows the architecture of Hu-Fu, which consists of three modules: *query rewriter*, *drivers*, and *query interface*. From a functional perspective, the query rewriter and drivers optimize the query *efficiency*, and the query interface improves the *usability* of Hu-Fu.

Query Rewriter (Sec. 4). It decomposes mainstream spatial queries (federated range query, range counting, kNN query, distance join, and kNN join) into plaintext operators (executed within silos) and secure operators

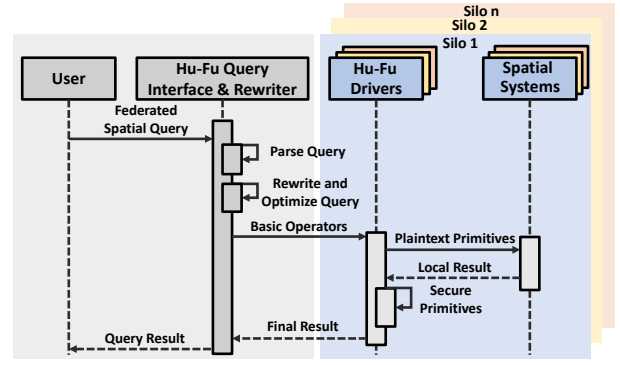


Fig. 3: Illustration of Hu-Fu workflow.

(executed across silos). We define three *plaintext* operators (plaintext range query, range counting, and kNN query) and five *secure* operators (secure summation, count comparison, set union, distance comparison, and location perturbation) as the basic operators. We also design novel execution plans that decompose these federated spatial queries into basic operators.

Drivers (Sec. 5). Hu-Fu’s drivers implement the basic operators defined in the query rewriter as efficient *primitives* that can adapt to heterogeneous spatial databases at the backend. Each operator is implemented by a specific primitive. Specifically, secure operators are implemented as *secure primitives* with dedicated optimizations [18, 23, 35, 37]. Plaintext operators are implemented as *plaintext primitives* on top of the underlying spatial databases, which support various systems, *e.g.*, PostGIS [10], SpatiaLite [11], MySQL [8], GeoMesa [27], Simba [64] and SpatialHadoop [21].

Query Interface (Sec. 6). This module (1) provides a transparent and unified federation view to users, and (2) supports both asymmetric and symmetric federated spatial queries written in SQL. We implement the query interface by extending the schema manager and parser of Calcite [19]. We also provide interfaces such as JDBC for easy integration of Hu-Fu to users’ programs.

3.2 Workflow

Fig. 3 depicts the workflow of Hu-Fu for querying a data federation of n silos. The *query interface* and *query rewriter* are deployed on the user’s machine to provide a portal for spatial services, while each silo runs a *driver* that interacts with its underlying spatial database.

In the workflow, a user’s SQL-based federated spatial query is first parsed by the *query interface*. Then, the *query rewriter* transforms and optimizes the query into a sequence of plaintext and secure operators. These operators are executed as primitives by the *drivers*:

plaintext primitives run locally on spatial databases to produce intermediate results, while secure primitives assemble these partial results to obtain the final answer, which is returned to the user via the query interface.

4 Query Rewriter

This section presents the design of the query rewriter in Hu-Fu, which decomposes a federated spatial query into basic operators to form the query execution plan. Specifically, we first introduce the involved basic operators in Sec. 4.1. Next, we explain the overall decomposition strategies in Sec. 4.2. Then, we introduce the rewriter of *asymmetric* federated spatial queries in Sec. 4.3 from two categories: *radius-known* and *radius-unknown* queries. Since asymmetric federated spatial queries do not assume query privacy, we also propose efficient solutions to *symmetric* federated spatial queries in Sec. 4.4. Finally, we discuss practical issues in Sec. 4.5.

4.1 Basic Operators

Our acceleration strategy is to *decompose queries into basic operators so that the majority of distance-related operations occur within silos in plaintext*, thereby reducing the need for secure operations across silos. The selection principle of basic operators is explained below.

4.1.1 Operator Selection Principles

There are two categories of basic operators in Hu-Fu: *plaintext* and *secure* operators. The *plaintext* operators handle local queries within each individual silo, while the *secure* operators perform atomic computations over sensitive data in a privacy-preserving manner.

- **Plaintext Operators.** They can involve the distance-related operations compulsory in spatial queries, but should be common operations widely supported by diverse spatial databases.
- **Secure Operators.** They should avoid distance-related operations unless strictly necessary, and efficiently implemented operators are preferable.

Adhering to these principles, we select three plaintext operators and five secure operators, which will be elaborated in Sec. 4.1.2 and Sec. 4.1.3, respectively.

4.1.2 Plaintext Operators

We define three plaintext operators: *plaintext range query*, *range counting*, and *kNN query*. These operators are performed within each silo F_i . Hence, they can be conducted in plaintext without compromising security.

Definition 1 (Plaintext Range Query/Counting) Given a silo F_i and a query range \mathcal{R} , the *plaintext range query* $\text{RQ}(F_i, \mathcal{R})$ retrieves the spatial objects in F_i that fall within \mathcal{R} , and the *plaintext range counting* $\text{RC}(F_i, \mathcal{R})$ returns the number of such objects.

Definition 2 (Plaintext kNN Query) Given a silo F_i , a query object q , and a positive integer k , the *plaintext kNN query* $\text{kNN}(F_i, q, k)$ retrieves the k nearest spatial objects in F_i to the query object q .

The plaintext operators comply with the principles described in Sec. 4.1.1, because they are supported by almost all spatial databases. They are implemented as *plaintext primitives* in Hu-Fu drivers, which we defer to Sec. 5.1. The query range can be various shapes, such as circles and rectangles. For ease of presentation, we mainly focus on circular ranges in this section and discuss extensions to other shapes in Sec. 4.5.

4.1.3 Secure Operators

We define five secure operators: *summation*, *count comparison*, *set union*, *distance comparison*, and *location perturbation*. The first three secure operators are designed to preserve data privacy, while the latter two secure operators aim to protect query privacy.

Definition 3 (Secure Summation) Given n data silos $\{F_i\}$ each holding a private value v_i , this operator SUM sums up these values, *i.e.*, $\text{SUM}(v_1, \dots, v_n) = \sum_{i=1}^n v_i$, while protecting the privacy of v_i in silo F_i from all other silos F_j ($\forall j \neq i$).

Definition 4 (Secure Count Comparison) Given n data silos $\{F_i\}$ each holding a private count v_i and a public constant k , this operator CMP compares the sum of these counts with k , *i.e.*, $\text{CMP}(v_1, \dots, v_n, k) = \text{sign}(\sum_{i=1}^n v_i - k)$, without leaking the sum $\sum_{i=1}^n v_i$ or the count v_i in silo F_i to any other silos F_j ($\forall j \neq i$).

Definition 5 (Secure Set Union) Given n data silos $\{F_i\}$ each holding a set of spatial objects $S_i = \{o_1^i, \dots, o_{m_i}^i\}$, this operator SUN computes the union of spatial objects from all silos, *i.e.*, $\text{SUN}(S_1, \dots, S_n) = \cup_{i=1}^n S_i$, without leaking the spatial objects S_i in silo F_i to any other silos F_j ($\forall j \neq i$).

Definition 6 (Secure Distance Comparison) Given a query user holding a private location l_q , a data silo holding a private location l_o , and a distance threshold r , this operator DCMP compares the distance between l_q and l_o with the threshold r , *i.e.*, $\text{DCMP}(l_q, l_o, r) = \text{sign}(\text{dist}(l_q, l_o) - r)$, without leaking the location of either party (*i.e.*, the user or silo) to the other.

Definition 7 (Secure Location Perturbation) Given a private location $x \in \mathbb{R}^2$, this operator Geol obfuscates it into a location $z = \text{Geol}(x)$ while satisfying (ϵ, δ) -Geo-Indistinguishability (Geo-I) [60]. The privacy requirement of (ϵ, δ) -Geo-I is satisfied iff the following inequality holds for any two locations x, x' in the location set and any location subset Z :

$$\Pr[\text{Geol}(x) \in Z] \leq e^{\epsilon \cdot \text{dist}(x, x')} \cdot \Pr[\text{Geol}(x') \in Z] + \delta \quad (1)$$

where $\Pr[\text{Geol}(x) \in Z]$ is the probability that the perturbed location belongs to the subset Z , and ϵ, δ represent the privacy preservation level of Geo-I.

ϵ -Geo-I [13] adapts the de facto standard privacy notion, *differential privacy* [39], to protect location data, where ϵ is known as the privacy budget. (ϵ, δ) -Geo-I relaxes the definition of ϵ -Geo-I by allowing a small failure probability δ . This way of relaxation has gained widespread usages in (standard) differential privacy [39].

Remark. The secure operators comply with the principles in Sec. 4.1.1 since (1) most of them do not involve distance operations (with the exception of secure distance comparison) and (2) all of them have dedicated and efficient implementations (see Sec. 5.2 for details).

4.2 Overview of Our Decomposition Strategies

In the following, we formally define the federated spatial queries and introduce our taxonomy to categorize them (Sec. 4.2.1). We then elaborate on the main ideas of our decomposition strategies for each category (Sec. 4.2.2).

4.2.1 Federated Spatial Queries and Taxonomy

Before diving into our decomposition strategies, we first define the five federated spatial queries. The privacy requirement below includes either data privacy alone or both data and query privacy defined in Sec. 2.1.

Definition 8 (Federated Range Query/Counting)

Given a federation F of n data silos $\{F_i\}$, and a query range \mathcal{R} , a *federated range query* retrieves all spatial objects located within \mathcal{R} , while a *federated range counting* returns the number of such objects. Both queries need to satisfy the privacy requirement.

Definition 9 (Federated Distance Join)

Given a federation F of n data silos $\{F_i\}$, an input dataset Q of spatial objects, and a distance radius r , a federated distance join retrieves all pairs of objects (q, o) where $q \in Q, o \in F$ such that the distance $\text{dist}(l_q, l_o) \leq r$, while satisfying the privacy requirement, i.e.,

$$Q \bowtie_r F = \{(q, o) \mid q \in Q, o \in F, \text{dist}(l_q, l_o) \leq r\}.$$

Definition 10 (Federated kNN Query/Join) Given a federation F of n data silos $\{F_i\}$, a query object q , and a positive integer k , a *federated kNN query* retrieves the k nearest objects in F to the query object q , i.e.,

$$\forall o \in \text{kNN}(q), \forall o' \in F \setminus \text{kNN}(q), \text{dist}(q, o) \leq \text{dist}(q, o').$$

When the query objects form an input dataset Q , a *federated kNN join* retrieves all pairs of objects (q, o) where $q \in Q$ and o belongs to the kNN of q in F , i.e.,

$$Q \bowtie_{\text{kNN}} F = \{(q, o) \mid q \in Q, o \in \text{kNN}(q)\}.$$

Both queries need to satisfy the privacy requirement.

Taxonomy. The above queries can be categorized from *two orthogonal dimensions*: the scope of the privacy requirement and whether the searching radius (of the query range) is explicitly given. Specifically, based on whether query privacy is included in the privacy requirement, the queries are classified into *asymmetric* and *symmetric* queries (see the differences in Sec. 2.1). Based on whether the searching radius is explicitly given, the queries are classified into *radius-known* and *radius-unknown* queries. Intuitively, the federated range query, range counting, and distance join belong to radius-known queries, while the federated kNN query and kNN join belong to radius-unknown queries.

4.2.2 Main Idea of Our Decomposition Strategies

Basic Principle. In Hu-Fu, the core principle of the query rewriter is to decompose federated spatial queries into as many plaintext operators and as few secure operators as possible such that a large portion of the query can be executed in plaintext without compromising security. At a high level, a federated spatial query is initially processed using plaintext operators within each silo, and their results are then securely assembled to form the final outcome. At the minimum, one secure operator is compulsory, and additional secure operators may be required if there are interactions across silos.

Based on the aforementioned basic operators and taxonomy of queries, we now introduce the *main ideas* of decomposing different categories of queries.

Main Idea for Asymmetric Queries with Data Privacy Solely. Our idea is elaborated as follows:

- **Radius-Known Queries.** A radius-known query (e.g., federated range query and range counting) can be decomposed into the corresponding plaintext operators within each silo and only one secure operator (e.g., a secure set union or a secure summation) for assembling the partial results across silos.

- **Radius-Unknown Queries.** Each radius-unknown query (*e.g.*, federated kNN query) is viewed as an iterative process of trialing different search radii until exactly k spatial objects are found within the circular search area. This execution can be converted into multiple radius-known queries (*e.g.*, federated range counting), with the number of radius-known queries minimized by a binary search. Each iteration utilizes a secure operator to ensure data privacy.

Key Insight for Symmetric Queries with Both Data and Query Privacy. When additionally considering query privacy, our key insights are as follows:

- **Radius-Known Queries.** A native solution employs a secure distance comparison operator for every spatial object, but leads to excessive secure distance operations. Instead, we first obfuscate the sensitive query location using a secure location perturbation operator to create a noised location that can be safely published to each silo. Then, leveraging the previous idea for radius-known queries, each silo identifies a small set of candidates. For each candidate, a secure distance comparison operator verifies if its distance to the query location is within the specified radius, while protecting query privacy.
- **Radius-Unknown Queries.** Similar to the aforementioned ideas for decomposing radius-unknown queries, we can still decompose them into a series of radius-known queries.

4.3 Decomposing Asymmetric Queries with Data Privacy Only

This subsection proposes our methods for decomposing radius-known queries (Sec. 4.3.1) and radius-unknown queries (Sec. 4.3.2), which only consider data privacy. The decomposition plans are summarized in Table 2.

4.3.1 Decomposing Radius-Known Queries

Among the five queries, the federated range query, range counting, and distance join are radius-known queries.

Decomposition Plan. (1) **Federated range query** can be decomposed into n *plaintext range queries*, each with a radius r , where each plaintext operator retrieves the partial result within each one of the n silos. Afterwards, a *secure set union* operator assembles these partial result while maintaining data privacy. (2) Similarly, **federated range counting** can be decomposed into n *plaintext range counting* operators to obtain n partial counts. These partial counts will later be aggregated by a *secure summation* operator. (3) **Federated distance**

join is equivalent to requesting federated range queries $|R|$ times, each of which follows the previous plan.

Complexity Analysis. Let T_{RQ} and T_{RC} denote the time complexity of plaintext range query and range counting, respectively. $|S|$ denotes the size of returned set. Based on the complexities of secure operators (see Sec. 5.2), the time complexity and communication cost of the radius-known queries are as follows. (1) *Federated range query* takes $O(T_{RQ} + n + |S|)$ time and $O(n + |S|)$ communication cost. (2) *Federated range counting* takes $O(T_{RC} + n^3)$ time and $O(n^2)$ communication cost. (3) *Federated distance join* takes $O(|R| \cdot T_{RQ} + n + |S|)$ time and $O(n + |S|)$ communication cost.

4.3.2 Decomposing Radius-Unknown Queries

Federated kNN query and kNN join are classified as radius-unknown queries due to the absence of an explicitly given range. Their decomposition plan is to first get an appropriate range and then filter the points in the range, as explained in detail below.

Decomposition Plan. Similar to the relation between federated range query and federated distance join in Sec. 4.3.1, federated kNN join can be viewed as $|R|$ independent federated kNN queries. Hence, we mainly explain how to decompose a federated kNN query.

- **Basic Idea.** Recall from Sec. 4.2, the strategy to decompose radius-unknown queries is to convert them into multiple rounds of radius-known queries. We first derive a radius r via a binary search and then retrieve the spatial objects within this search range. For each radius r , we securely check whether the counting result is smaller than k . As long as r falls between the k th and the $(k + 1)$ th nearest distance to the query object q , the spatial objects within this range are precisely the k nearest neighbors.
- **Algorithm Details.** Alg. 1 illustrates the decomposition of a federated kNN query. Lines 1-8 derive the radius r . We initialize a lower bound ($l = 0$) and upper bound ($u = U$) of the radius, where U can be set as the spatial area's diameter or a user-defined value. A binary search is then performed to find the appropriate radius until reaching the distance precision ϵ_0 (lines 2-9). In each iteration, r is set to $(l + u)/2$. For each r , a *plaintext range counting* operator is executed within each silo, and a *secure count comparison* operator is invoked to compare the total count with the integer k (lines 4-5). If the total count is less than k (*i.e.*, $sign < 0$), indicating an undersized radius, l will be increased to r . Conversely, if the total count exceeds k (*i.e.*, $sign > 0$), u will be decreased to r . The binary search ensures that the final radius r is sufficiently close to the k th

Table 2: The number of basic operators in the decomposition plans of *asymmetric* federated spatial queries. Radius-known queries only involve one type of secure operators (secure summation or set union). Radius-unknown queries are executed in multiple rounds which additionally require secure count comparisons to ensure security.

Category	Federated Spatial Query	#(Plaintext Operator)		#(Secure Operator)	
		Range Query	Range Counting	Count Comparison	Set Union/Summation
Radius-Known	Range Query	n	0	0	1/0
	Range Counting	0	n	0	0/1
	Distance Join	$n R $	0	0	$ R /0$
Radius-Unknown	kNN Query	n	$O(n \log U)$	$O(\log U)$	1/0
	kNN Join	$n R $	$O(n R \log U)$	$O(R \log U)$	$ R /0$

n is the number of silos, R is the input dataset in spatial joins, and U is the upper bound for the binary-search radius.

Algorithm 1: Asymmetric federated kNN query

Input: federation F , query object q , integer k
Output: the (exact) query answer ans

```

1  $[l, u] \leftarrow [0, U]$ , where  $U$  is a predefined upper bound;
2 while  $u - l \geq \epsilon_0$  do
3    $r \leftarrow (l + u)/2$ ,  $\mathcal{R} \leftarrow \text{circle}(q, r)$ ;
4   foreach  $\text{silos } F_i \in F$  do // perform in parallel
5      $v_i \leftarrow \text{plaintext range counting } \text{RC}(F_i, \mathcal{R})$ ;
6    $\text{sign} \leftarrow \text{secure count comparison } \text{CMP}(\{v_i\}, k)$ ;
7   if  $\text{sign} < 0$  then  $l \leftarrow r$ ;
8   else if  $\text{sign} > 0$  then  $u \leftarrow r$ ;
9   else break;
10 foreach  $\text{silos } F_i \in F$  do // perform in parallel
11    $S_i \leftarrow \text{plaintext range query } \text{RQ}(F_i, \text{circle}(q, r))$ ;
12 return  $ans \leftarrow \text{secure set union } \text{SUN}(S_1, \dots, S_n)$ ;
```

nearest distance. Finally, a *plaintext range query* is executed on each silo, and the partial results are collected using a *secure set union* (lines 9-10).

In Alg. 1, the distance precision ϵ_0 is initially set based on the application requirement. For example, many spatial applications (*e.g.*, taxi-calling) have a minimum distance precision requirement that is typically measured in meters. Then, ϵ_0 can be set to 1 meter.

Complexity Analysis. Alg. 1 requires $O(\log \frac{U}{\epsilon_0}) = O(\log U)$ iterations to obtain the final radius, where ϵ_0 is a constant to denote this radius's precision. In each iteration, the plaintext range counting takes $O(T_{\text{RC}})$ time, and the secure count comparison takes $O(n)$ time and $O(n^2)$ communication cost. In lines 9-10, Alg. 1 performs a plaintext range query that takes $O(T_{\text{RQ}})$ time and a secure set union that takes $O(n + k)$ time and $O(n + k)$ communication cost. Overall, the total time complexity is $O((T_{\text{RC}} + n) \cdot \log U + T_{\text{RQ}} + k)$, and the communication cost is $O(n^2 \cdot \log U + k)$. Intuitively, the complexity of federated kNN join is equal to that of federated kNN query, multiplied by a factor $|R|$.

Example 3 Fig. 4 illustrates the procedure of Alg. 1 with a query point $(4, 4)$ and $k = 3$ over 3 silos, where the objects with the same color belong to the same silo.

The query rewriter decomposes this query into multiple rounds of radius-known queries. In the 1st round, a plaintext range counting with center $(4, 4)$ and radius 4 is sent to each silo and a secure count comparison with k is performed across silos. And we get 10 objects, which is greater than k . Hence in the 2nd round, the radius decreases to 2 and is sent to each silo for plaintext range counting and secure count comparison. There are 2 objects, which is fewer than k . Thus, in the 3rd round, the radius increases to 3, and the procedure continues, where the secure count comparison results implies that $\text{sign} = 0$ and the search terminates. Finally, the basic operators, including the plaintext range query with the center $(4, 4)$ and radius 3 and secure set union, are performed to retrieve the 3 query answers.

Optimization via Differential Privacy. We exploit differential privacy [39] to further accelerate federated kNN query and federated kNN join from two aspects.

- **Tighten Predefined Upper Bound.** We ask each F_i to perform a plaintext kNN query operator and return the k th object's distance U_i to the query point. Since directly returning such values may violate the data privacy requirement, we apply the truncated Laplacian mechanism [15] on it. That is, let each silo add a positive noise and obtain the perturbed value \tilde{U}_i . We can tighten the upper bound as the shortest distance in all silos, *i.e.*, $U = \min\{\tilde{U}_i\}$, since there must be at least k objects in this range.
- **Reduce Running Time and Communication Cost in Secure Count Comparison.** The secure count comparison in Alg. 1 compares $\sum_1^n v_i$ with k , resulting in $O(n^2)$ running time and communication cost. However, when $\sum_1^n v_i$ differs significantly from k , this can be reduced to $O(n)$ by using the Laplacian mechanism [39] in differential privacy. This mechanism injects a noise into the local count in each silo, and then perturbed counts are aggregated in plaintext. If the perturbed total count

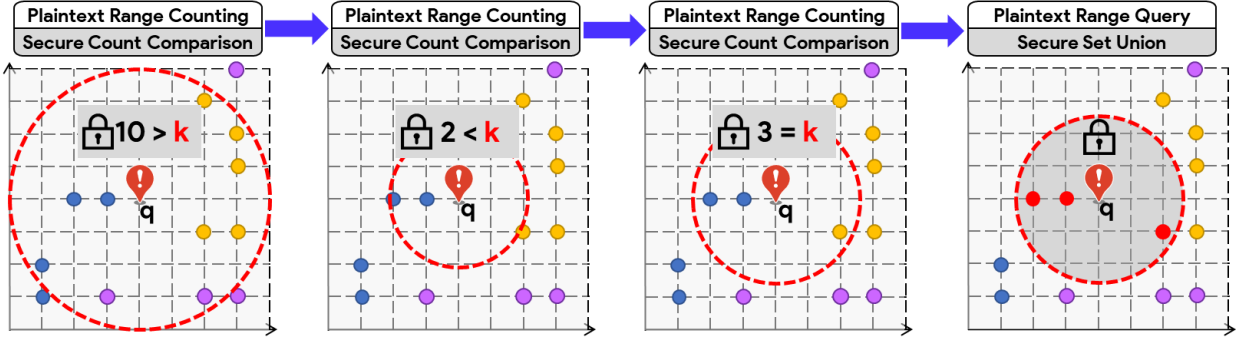


Fig. 4: Running example for (asymmetric) federated kNN query ($k = 3$).

Algorithm 2: Symmetric federated range query

Input: federation F , query object q , radius r
Output: the (exact) query answer ans

- 1 $q' \leftarrow$ secure location perturbation $\text{Geol}(q)$;
- 2 $r' \leftarrow r + \text{dist}(l_q, l_{q'})$, range $\mathcal{R}' \leftarrow \text{circle}(q', r')$;
- 3 **foreach** silo $F_i \in F$ **do** // perform in parallel
- 4 $Cand_i \leftarrow$ plaintext range query $\text{RQ}(F_i, \mathcal{R}')$
- 5 **for** silo $F_i \in F$ **do** // perform in parallel
- 6 **foreach** candidate spatial object $o \in Cand_i$ **do**
- 7 $sign \leftarrow$ secure distance comparison $\text{DCMP}(l_q, l_o, r)$;
- 8 **if** $sign \leq 0$ **then** $ans \leftarrow ans \cup \{o\}$;

is much smaller or larger than k , we directly adjust the threshold without running the secure operator.

4.4 Decomposing Symmetric Queries with both Data Privacy and Query Privacy

This subsection presents our methods for decomposing radius-known queries (Sec. 4.4.1) and radius-unknown queries (Sec. 4.4.2) with both data and query privacy. The decomposition plans are summarized in Table 3.

4.4.1 Decomposing Radius-Known Queries

Since query location must be protected in symmetric queries, radius-known queries can no longer be decomposed into plaintext range query/counting within each silo directly. Instead, we use the Geo-I mechanism [13, 60] to preserve the query privacy, as detailed below.

Decomposition Plan for Federated Range Query.

Alg. 2 presents the decomposition plan for federated range queries. Initially, a *secure location perturbation* operator is applied to generate an obfuscated object q' . Next, the search radius is increased by the distance from location l_q to $l_{q'}$ (line 2). This ensures that the expanded query range, denoted by a circle centered at q' with a radius $r' = r + \text{dist}(l_q, l_{q'})$, completely covers the intended query area. Within each silo, a *plain-*

text range query is then performed using the expanded query range (line 3). As a result, each silo obtains a set of candidates for the query answer. To refine these candidates, *secure distance comparison* operators are employed to filter out those outside the true query range and collect the final answer while satisfying both data privacy and query privacy (lines 4-7).

Example 4 Fig. 5a presents an illustrative example for Alg. 2. Suppose the query object q is located at $(2.5, 2.5)$ and the radius r of the circular query range is 1.6. By using the *secure location perturbation* operator, q is obfuscated into q' located at $(4, 4)$, so the radius r' is increased to $1.6 + \sqrt{(2.5 - 4)^2 + (2.5 - 4)^2} = 3.7$ (lines 1-2 of Alg. 2). After executing the *plaintext range query* operator with the expanded query range, we identify 3, 2, and 5 candidates (marked in different colors) in all three silos (line 3). Finally, each candidate is further refined by the *secure distance comparison* operator.

Extension to Other Radius-Known Queries. The decomposition plan for *federated range counting* is almost identical to Alg. 2. The key difference lies in line 7, where federated range counting only needs to aggregate the counts. When dealing with *federated distance join*, it is initially converted into a series of (symmetric) *federated range queries*. Subsequently, each federated range query is decomposed by Alg. 2.

4.4.2 Decomposing Radius-Unknown Queries

Similar to the binary search procedure in Alg. 1, a symmetric *federated kNN query* can be broken down into multiple rounds of (symmetric) radius-known queries. Besides, *federated kNN join* can still be decomposed into a series of independent federated kNN queries. Thus, we focus primarily on the necessary modifications for federated kNN queries in the following.

Naive Decomposition Plan. A naive extension of Alg. 1 can be time-consuming due to the trivial upper

Table 3: The number of basic operators in the decomposition plans of *symmetric* federated spatial queries.

Category	Federated Spatial Query	#(Plaintext Operator)		#(Secure Operator)		
		Range Query	kNN Query	Location Perturbation	Distance Comparison	Set Union
Radius-Known	Range Query	n	0	1	$ Cand $	1
	Range Counting	n	0	1	$ Cand $	0
	Distance Join	$n R $	0	$ R $	$ R \cdot Cand $	$ R $
Radius-Unknown	kNN Query	$O(n \log U)$	n	$1 + n$	$O(Cand \log U)$	1
	kNN Join	$O(n R \log U)$	$n R $	$(1 + n) R $	$O(Cand \log U \cdot R)$	$ R $

n is the number of silos, $|R|$ is the size of the input dataset R in spatial joins, U is the upper bound for the binary-search radius, and $|Cand|$ is the total number of candidate objects from all data silos.

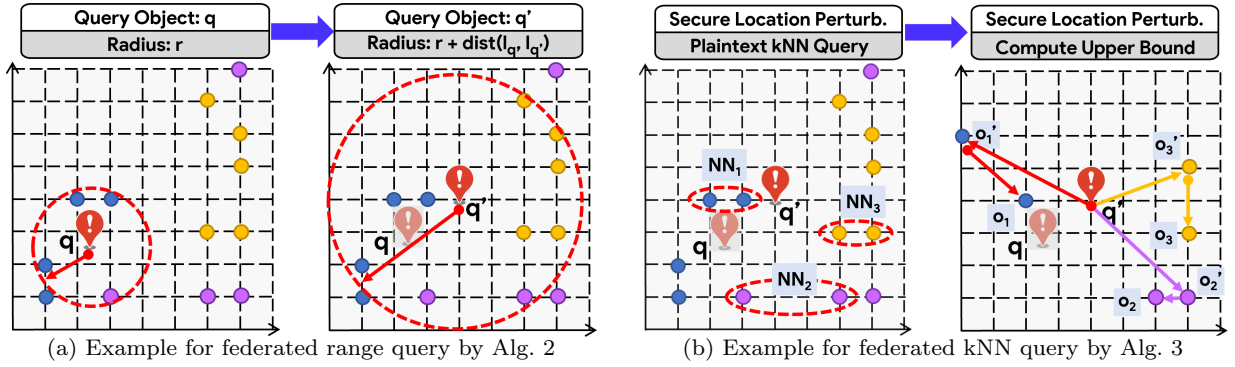


Fig. 5: Examples for decomposing (symmetric) federated range query and kNN query

bound of the search radius. If the initial upper bound is set too high, the decomposed radius-known queries will require a large number of secure distance operations, which becomes the major efficiency bottleneck.

Optimized Decomposition Plan. To overcome the limitation of the naive method, we devise Alg. 3 to compute a tighter upper bound based on the perturbed location. Specifically, line 1 perturbs the query object q into q' privately. In line 2, plaintext kNN query is performed in each silo to identify the k nearest neighbors to q' . However, we cannot send the k th nearest distances to the user's client as the upper bound, since it leaks information about locations in the data silos. Instead, each silo obfuscates the k th nearest neighbor $o_i \in NN_i$ into a noised spatial object o'_i (lines 3-4). In line 5, each silo locally computes its own upper bound U_i and sends it to the user's client. Finally, the upper bound can be safely set as the minimum value among $\{U_i\}$ (line 6).

Example 5 Fig. 5b illustrates Alg. 3. In line 1, the query object q located at (2.5, 2.5) is perturbed into the object q' located at (4, 4). Each silo then performs a *plaintext kNN query* with $k = 2$ on q' . The top-2 nearest neighbors in three silos are denoted by NN_1 - NN_3 , where o_i is the farthest object to q' within NN_i . For example, the object o_1 at (2, 4) has a distance of 2 to q' . However, revealing this distance directly may leak spatial

Algorithm 3: Compute tight upper bound for optimizing symmetric federated kNN query

Input: federation F , query object q , integer k
Output: the upper bound U for binary-search radius
1 $q' \leftarrow$ secure location perturbation $\text{Geol}(q)$;
2 **foreach** silo $F_i \in F$ **do** // perform in parallel
3 $NN_i \leftarrow$ plaintext kNN query $\text{kNN}(F_i, q', k)$
4 $o_i \leftarrow \arg \max_{o \in NN_i} \{\text{dist}(l_o, l_{q'})\}$;
5 $o'_i \leftarrow$ secure location perturbation $\text{Geol}(o_i)$;
6 $U_i \leftarrow \text{dist}(l_{q'}, l_{o'_i}) + \text{dist}(l_{o'_i}, l_{o_i})$;
7 **return** $U \leftarrow \text{dist}(l_q, l_{q'}) + \min\{U_i \mid i = 1, \dots, n\}$;

information about o_1 . Instead, Alg. 3 leverages a *secure location perturbation* operator to obfuscate o_1 to o'_1 located at (0, 6). Similarly, o_2 and o_3 are perturbed to o'_2 and o'_3 , respectively. Based on the Euclidean distances, we have $U_1 = 7.30$, $U_2 = 5.16$, and $U_3 = 5.24$ (line 5). Finally, we pick the minimum from $\{U_i\}$ and derive the tight upper bound $U = 5.16 + 2.12 = 7.28$.

The correctness of Alg. 3 is proved in Lemma 1.

Lemma 1 The upper bound U in Alg. 3 is no shorter than the k th nearest distance to the query object q in the data federation F .

Proof Let d^* denote the k th nearest distance to q , so d^* should satisfy the following inequality:

$$d^* \leq \max\{\text{dist}(l_q, l_o) \mid o \in NN_i\}, \forall i \in [1, n] \quad (2)$$

According to the triangle inequality, we have

$$\text{dist}(l_q, l_o) \leq \text{dist}(l_q, l_{q'}) + \text{dist}(l_{q'}, l_o) \quad (3)$$

Based on Eq. (2) and Eq. (3), we can derive that

$$d^* \leq \text{dist}(l_q, l_{q'}) + \max\{\text{dist}(l_{q'}, l_o) \mid o \in NN_i\}, \forall i \in [1, n]$$

Since $\max\{\text{dist}(l_{q'}, l_o) \mid o \in NN_i\} = \text{dist}(l_{q'}, l_{o_i})$ based on the line 3 of Alg. 3, we have

$$d^* \leq \text{dist}(l_q, l_{q'}) + \text{dist}(l_{q'}, l_{o_i}), \forall i \in [1, n]$$

Based on the triangle inequality for $\text{dist}(l_{q'}, l_{o_i})$ and the definition of U_i in line 5 of Alg. 3, we have

$$\begin{aligned} d^* &\leq \text{dist}(l_q, l_{q'}) + (\text{dist}(l_{q'}, l_{o'_i}) + \text{dist}(l_{o'_i}, l_{o_i})) \\ &\leq \text{dist}(l_q, l_{q'}) + U_i \end{aligned} \quad (4)$$

According to the inequality in Eq. (4) and the definition of U in line 6, we can now prove $d^* \leq U$.

Remark. In Alg. 3, the k nearest neighbors NN_i of the perturbed location q' are used to derive the upper bound. Notice that $\{NN_i\}$ do not necessarily encompass all the query results of the original location q . By contrast, with Lemma 1, the federated range query/counting during the binary-search can ensure that the candidate set includes the exact kNN of q .

4.5 Discussion

We provide further discussions on the query rewriter.

Security of Query Rewriter. We prove the security of our query rewriter based on the composition lemma in [30]. The idea is to show the decomposition plans for radius-known queries and radius-unknown queries will not reveal any extra information other than the final result due to the usage of secure operators. We also present a case study that proves it is hard for a semi-honest adversary to attack Hu-Fu. Please refer to Appendix B for the detailed security proof and case study of protecting security under specific semi-honest attack.

Handling Ties in kNN Queries. In federated kNN queries, we may encounter ties where multiple spatial objects share the same distance (*i.e.*, the k th nearest distance r^*) to the query object q . Here, we must resolve **two technical issues**: (1) identifying the presence of ties, and (2) retrieving exactly k nearest neighbors.

(1) Line 8 ($\text{sign} = 0$) of Alg. 1 indicates the current search radius covers exactly k objects (*i.e.*, no ties). If $\text{sign} \neq 0$ during the binary-search, then there are ties.

(2) Once ties are identified, we proceed to retrieve exactly k spatial objects. Let $[l, u]$ denote the lower and

upper bounds of r^* . First, we use a federated range query with the circular range circle(q, l) to cache the nearest neighbors that are not part of the ties. We denote the number of these objects as k_l . Next, we select $k - k_l$ objects from the tied ones by sequentially requesting objects from all data silos until we reach the desired count. Finally, we use a secure union operator to collect the cached partial answers from all data silos.

Please refer to Appendix C for the complete pseudocode, example illustrations, and experimental evaluation of our solution to handle the case of ties.

Extension to Rectangular Query Range. The decomposition plan for radius-known queries can be extended to accommodate rectangular-shaped of query ranges. For *asymmetric queries*, the extension can be seamlessly implemented by using plaintext range query/counting operators for rectangular query ranges, which are typically supported by spatial database systems. For *symmetric queries*, where the query object (*i.e.*, the rectangle center) is private, our extension proceeds as follows. We first compute the rectangle's circumscribed circle. Next, by querying the circumscribed circle with lines 1-4 of Alg. 2, we identify potential candidates. Finally, we securely verify if a candidate (x, y) lies within the rectangle $[(x_L, x_R), (y_L, y_R)]$ using the Yao's garbled circuit (GC) protocol [24] to check the inequalities $x \geq x_L$, $x \leq x_R$, $y \geq y_L$, and $y \leq y_R$. The Yao's GC protocol here can be implemented using OblivM [42].

Beyond Mainstream Spatial Queries. The query rewriter also supports *aggregation queries*, *e.g.*, the aggregate attribute on the result of kNN query or range query. For example, the range aggregate query can be decomposed similarly to a federated range counting. Our solution can be also extended to support *approximate spatial queries* by replacing the exact methods for the plaintext operators with approximate ones. While this is easy to implement, it may be hard to achieve a good balance between efficiency and accuracy.

5 Drivers

In Hu-Fu, a driver is deployed on each data silo, consisting of both *plaintext primitives* (Sec. 5.1) and *secure primitives* (Sec. 5.2). Here, plaintext primitives refer to the implementations of plaintext operators that leverage the local spatial database at each silo. Secure primitives, on the other hand, indicate our secure protocols tailored for the secure operators defined in Sec. 4.1.3.

Unlike existing systems [14, 56], Hu-Fu aims to support **heterogeneous** databases through drivers. In this way, Hu-Fu can *enhance usability* and *avoid costly data migration* compared to these solutions that assume local databases are **homogeneous**. To achieve this, the

main difficulties include (1) drivers must integrate with the query rewriter to convert plaintext operators into diverse query formats used by data silos, and (2) drivers must offer default implementations of plaintext operators for local databases that lack support.

5.1 Plaintext Primitives

Plaintext primitives implement *plaintext range query*, *range counting*, and *kNN query*. They are implemented as an interface on top of the underlying spatial databases for portability and to harness existing range query and range counting implementations.

Primitive Implementation. The plaintext primitives are implemented by the underlying spatial databases.

- For databases that support these plaintext queries, *e.g.*, Simba [64] and PostGIS [10], we utilize the built-in functions for these queries or generate the corresponding SQL request. For example, in PostGIS [10], a plaintext range counting on silo F_i with the center p and radius r of a circular range can be implemented by requesting the SQL below.

```
SELECT COUNT(*) FROM  $F_i$ 
WHERE ST_DWithin(p,  $F_i$ .location, r);
```

- When databases lack native support for any query, the drivers offer a default implementation based on their supported queries and indexes. For example, GeoMesa [27] does not inherently support range counting, so we extend range counting by calling a range query and subsequently counting the result size.

Time Complexity. In modern spatial databases, plaintext range query, range counting, and kNN query can take $O(\log m + |S|)$, $O(\log m)$, and $O(\log m)$ time [44], where m is the data size and $|S|$ is the output size.

Remark. In practice, the actual performance of plaintext primitives depends on the native implementation of the local spatial database at each silo. Thus, when silos utilize heterogeneous spatial databases, the efficiency of federated spatial queries can be limited by the slowest plaintext primitive (see Sec. 7.5).

5.2 Secure Primitives

The secure primitives, including secure summation, count comparison, set union, distance comparison, and location perturbation, are independent of local databases.

Primitive Implementation. Each secure primitive is optimized with a tailored secure protocol as follows.

Secure Summation. This primitive is based on [23]. Initially, each silo F_i holds a private value v_i and

all n silos agree on n distinct public parameters $\{u_i\}$. Each silo F_i then selects a random polynomial of degree $n-1$ in the form $t_i(x) = (\sum_{k=1}^{n-1} a_{ik}x^k) + v_i$, where a_{ik} is the random coefficient independently generated by silo F_i , and v_i denotes the private value (*i.e.*, local count) of silo F_i . These variables are kept secret from others by silo F_i . Next, each silo F_i evaluates its polynomial at the public parameters $\{u_1, \dots, u_n\}$ and sends the resulting value $t_i(u_j)$ to every other silo F_j . Once the silo F_j receives all values $\{t_i(u_j) | i \neq j\}$ from the other silos, we have $S(u_j) = \sum_{i=1}^n t_i(u_j) = (\sum_{k=1}^{n-1} (u_j)^k \sum_{i=1}^n a_{ik}) + \sum_{i=1}^n v_i$. Afterward, this silo sends $S(u_j)$ to the query user. The user can interpret each $S(u_j)$ as a linear equation $S(u_j) = \sum_{k=1}^{n-1} (u_j)^k z_k + z_n$ in n unknown variables z_k , where $z_k = \sum_{i=1}^n a_{ik}$ (for $k < n$) and $z_n = \sum_{i=1}^n v_i$. Now, the user can solve the system of linear equations using the received coefficients $\{u_j\}$ and constants $\{S(u_j)\}$ via Gauss elimination, and obtain the unknown variable z_n (*i.e.*, the sum of v_i).

Secure Count Comparison. The primitive compares the constant k with the sum of each silo F_i 's private range count v_i and prevents the leakage of either v_i or $\sum_{i=1}^n v_i$ to the silo F_j and the query user. The *main idea* is evaluating $X(\sum_{i=1}^n v_i - k)$ rather than directly computing $\sum_{i=1}^n v_i - k$ to avoid disclosing the actual sum of v_i , where X is a positive random number. Next, we implement this secure primitive by using existing secure multiplication protocol [18]. Specifically, this protocol [18] assumes that two multiplicands, X and Y , are partitioned into n shares x_i and y_i , where $X = \sum_{i=1}^n x_i$ and $Y = \sum_{i=1}^n y_i$. Each silo F_i holds the corresponding shares x_i and y_i , where x_i is randomly generated by this silo and $y_i = v_i - \frac{k}{n}$. Together, the multiplication XY happens to be $(\sum_{i=1}^n x_i)(\sum_{i=1}^n v_i - k)$. Finally, the comparison result is inferred from the sign of XY .

Secure Set Union. We implement this primitive based on the two-phase union method in [35] with additional optimizations. In the first phase, each silo appends its results into a global set, along with some fake records. Then, in the second phase, these fake records are removed from the set. To reduce the communication cost, the number of fake records should be as few as possible. Thus, we use the Laplace mechanism [39] in differential privacy to control the number of fake records. Moreover, by splitting the global set into batches, parallel executions are enabled for each silo to independently append and remove fake records from each batch, thereby resulting in a shorter latency.

Secure Distance Comparison. We leverage fully homomorphic encryption (FHE), the BGV scheme [12], to implement this primitive in three key steps.

(1) **Encrypt User's Data:** the query user encrypts their location (x_q, y_q) and the threshold r using the

public key, *i.e.*, $E(x_q), E(y_q), E(r)$, where $E(\cdot)$ is the encryption function. The encrypted data and public key are then sent to the data silo.

(2) **Compute Garbled Distance Difference:** by using FHE, the data silo computes the encrypted difference $E(\Delta)$ between the square of distance $\text{dist}(l_q, l_o)$ and square of threshold r as follows:

$$\Delta = \text{dist}(l_q, l_o)^2 - r^2 = (x_q - x_o)^2 + (y_q - y_o)^2 - r^2$$

$$E(\Delta) = (E(x_q) - E(x_o))^2 + (E(y_q) - E(y_o))^2 - E(r)^2$$

To further obfuscate the value of $E(\Delta)$, the silo applies a random polynomial function $f(\cdot)$ that only has odd powers and positive coefficients in each term. The obfuscation here prevents the user from inferring the exact distance $\text{dist}(l_q, l_o)$ after decryption, thereby protecting the silo's location privacy.

(3) **Decrypt:** upon receiving $f(E(\Delta))$ from the silo, the user decrypts it with the secret key and obtains $f(\Delta)$. The final result is derived based on the sign of $f(\Delta)$ without knowing the exact value of Δ .

Secure Location Perturbation. We implement this primitive based on the BPL mechanism in [60]. This mechanism obfuscates the original location (x, y) in the polar coordinate system. The polar angle θ is uniformly sampled from $[0, 2\pi]$. The polar radius r is sampled based on the -1 branch of the Lambert W function. If the sampled radius r is too long, it will be truncated into a random value in $[0, R]$, where R is a safe upper bound of radius based on the privacy parameters ϵ, δ . The resulting perturbed locations are $(x + r \cos \theta, y + r \sin \theta)$.

Complexity Analysis. The *secure summation* takes $O(n^3)$ time and $O(n^2)$ communication cost. The *secure count comparison* requires $O(n)$ time and $O(n^2)$ communication cost. The time complexity and communication cost of *secure set union* are $O(n + |S|)$, where $|S|$ is the output size. The *secure distance comparison* (between two parties) takes $O(1)$ time and communication cost, since the complexity of the BGV scheme [12] used for this primitive primarily depends on constant security parameters. The time complexity and communication cost of *secure location perturbation* are also $O(1)$ due to the usage of differential privacy mechanism.

6 Query Interface

For easy usability, the query interface of Hu-Fu offers a unified federation view to users (Sec. 6.1) and supports federated spatial queries in SQL (Sec. 6.2).

6.1 Unified Federation View

Hu-Fu's query interface provides a federation view to users, while the detailed information of silos is hidden. This not only enables users to send queries without caring about the silo organization, but also protects the data privacy of individual silos.

We implement this unified federation view by extending the schema manager of Calcite [19], a popular query processing framework. In Calcite's schema manager, each table is independent and indivisible. We treat silos as an abstraction layer below the table of schema manager. This means each table comprises multiple silo objects, and each object records the identity information of its silo. The silo identities are used when executing secure primitives. Specifically, the query rewriter will attach the identity information of all silo-level tables in the table of schema manager when distributing secure operators. Each silo only executes the corresponding secure primitives if the attached identity information matches the one locally stored.

6.2 Federated Spatial Queries in SQL

Based on the unified federation view, Hu-Fu query interface supports federated spatial queries in SQL by extending the SQL parser of Calcite with four keywords: `DWithin`, `kNN`, `Private_DWithin`, and `Private_kNN`. The first two keywords are used in asymmetric queries, and the last two are used in symmetric queries.

For example, an asymmetric federated range counting on a circular range centered at the point p with radius r can be expressed in SQL as

```
SELECT COUNT(*) FROM F
WHERE DWithin(p, F.location, r)
```

The `WHERE` clause checks whether the distance from p to an object in F is shorter than r . Similarly, an asymmetric federated kNN join on a relation R and federation F can be written in SQL as

```
SELECT R.id, F.id
FROM R JOIN F
ON kNN(R.location, F.location, k)
```

The `WHERE` clause indicates whether a spatial object in F belongs to the kNN set of the query point $o \in R$.

In contrast, when locations in both R and F need protection, a symmetric federated kNN join in SQL is

```
SELECT R.id, F.id
FROM R JOIN F
ON Private_kNN(R.location, F.location, k)
```

Other federated spatial queries can be written as SQL similarly with these four keywords.

7 Evaluation

In this section, we first introduce the experimental setup (Sec. 7.1), and then present the overall performances of asymmetric queries (Sec. 7.2) and symmetric queries (Sec. 7.3), scalability tests (Sec. 7.4), and results with heterogeneous spatial databases across silos (Sec. 7.5).

7.1 Experimental Setup

Datasets. Experiments are conducted on two datasets, with each object having a location and unique ID.

- **Multi-company Spatial Data in Beijing (BJ).** This dataset² was collected by 10 companies in Beijing, in June 2019, which has 1,029,081 spatial objects in total. The locations of these objects fall into an area from 39.5°N \sim 42.0°N and 115.5°E \sim 117.2°E. We use the dataset to simulate a real-world federation, where each company can be naturally regarded as a silo. During the evaluation, we vary the silo number n and queries without altering the spatial object distributions across silos.
- **OpenStreetMap (OSM).** This is a popular open dataset to evaluate spatial queries. We mainly use this dataset in the scalability test, where we sample 10^4 - 10^9 spatial objects from the Asia dataset in the OpenStreetMap [9]. Specifically, to simulate the spatial overlaps as in the BJ dataset, we assign a random silo ID for each point in the dataset and make each silo have the same number of data points.

Please refer to Appendix F for experiments on geographically partitioned dataset.

General-Purpose Baselines. As a data federation system, the evaluation first aims to compare Hu-Fu with existing general-purpose data federation systems: the GIS extensions of SMCQL [14] and Conclave [56].

- **SMCQL-GIS.** It adopts the principles of SMCQL [14], a garbled circuit (GC) based solution for relational data, to support spatial queries. We implement it with ObliVM [42], which is used in SMCQL for GC protocols across two silos (only) [56, 59]. Thus, it is only evaluated over two data silos.
- **Conclave-GIS.** It adopts the principles of Conclave [56], the secret sharing (SS) based solution for relational data, to support spatial queries. It is implemented with a different SS based library, MP-SPDZ [36], rather than Sharemind [18] in the original Conclave, since Sharemind is devised for only three silos [24] and it is a commercial library. In

contrast, MP-SPDZ is a popular open-source library that supports more than three silos based on SS.

- **SMCQL-GISext & Conclave-GISext** are their variants without assuming an honest broker, and uses our secure set union to assemble results.

These secure baselines implement federated spatial queries by exploiting similar queries for relational data in SMCQL or Conclave. Our extensions follow the strategy of having plaintext spatial queries within each silo's database and securely computing the final results. Specifically, for *federated range query*, these baselines execute plaintext range query in each silo and collect the partial results by either the honest broker or our secure set union. For *federated range counting*, they execute plaintext range counting and use secure summation to compute the final result. For *federated kNN query*, we regard it as a top-k query with a user-defined function (UDF). For example, each silo runs plaintext kNN query to compute k candidate neighbors along with their distances to the query object. Then, all n silos securely find the k nearest neighbors among nk candidates. For *federated distance join/kNN join*, we refer to their query plans for join queries and regard a federated distance/kNN join as multiple federated range/kNN queries.

Specialized Baselines. Beyond these general-purpose data federation systems, the evaluation also compares Hu-Fu with the following specialized baselines.

- **Additional Baselines for Asymmetric Queries.** The plaintext baseline Public directly collects local results from each silo without any secure operation, and serves as the upper bound of query efficiency.
- **Additional Baselines for Symmetric Queries.** We consider two more baselines for symmetric queries: LFHE [37] and PINED-RQ++ [48]. LFHE [37] is an industrial solution that utilizes Leveled Fully Homomorphic Encryption (LFHE) and two mutually untrusted servers to securely answer (exact) kNN queries over multiple data silos. This solution can be easily extended to support secure range query and counting over a spatial data federation. By contrast, PINED-RQ++ [48] leverages a differentially private index (*e.g.*, grid index for spatial data) and AES encryption [30] to approximately answer range queries with small errors. However, this method assumes that the user has access to data objects outside the query answer, potentially violating the data privacy requirement. Nevertheless, we select PINED-RQ++ [48] as a baseline for comparison, since it also utilizes differential privacy for filtering before verifying each candidate through encryption.

Metrics. We assess the query efficiency by two metrics:

² <https://share.weiyun.com/z4QfVhVv>

Table 4: Improvement with DP in federated kNN.

Silo number	2	4	6	8	10
Running time	2.9%	10.3%	19.6%	16.2%	14.0%
Communication	32.6%	31.5%	27.4%	39.4%	47.7%

(1) **Running time** is the time cost from receiving the query to returning the query answer to the user.

(2) **Communication cost** is the total network communication among the user and all data silos.

Implementation. We use PostgreSQL 10.15 with PostGIS extension as the default spatial database for all silos. To show the support of heterogeneous spatial data systems by Hu-Fu, we also use MySQL 5.7 [8], SpatialLite [11], GeoMesa 3.0.0 [27], Simba 1.0 [64], and SpatialHadoop 2.4.3 [21] as different silos, as will be explained in Sec. 7.5. They all use spatial indexes (R-Tree in PostGIS, Simba, SpatialHadoop, and MySQL, and R*-Tree in SpatialLite, and Z-Curve in GeoMesa) to speed up plaintext primitives by up to 2042 \times (see Appendix G). Among the compared solutions, LFHE and PINED-RQ++ are implemented in C++, while the others are implemented in Java. The reason for using C++ for LFHE and PINED-RQ++ is due to the lack of robust and open-source libraries in Java for the encryption methods (*e.g.*, CKKS [12]) they utilize.

7.2 Experiments on Asymmetric Queries

Parameter Setting. In this experiment, we compare the efficiency of different methods for all five federated spatial queries on the real dataset BJ. All the query points are randomly sampled from the dataset. We vary the number of silos from 2 to 10, and also test the impact of query-specific parameters. We set k to 16 for federated kNN query and kNN join, and the default query area of federated range query, range counting and distance join as 0.001%, and vary them from 4 to 64 and 0.00001% to 0.1% respectively. The range of these query-specific parameters is aligned with previous study [64]. When evaluating the query-specific parameters, we use 6 silos by default.

Environment. We run this experiment on a cluster of 11 machines. Each machine has 32 Intel(R) Xeon(R) Gold 5118 2.30GHz processors and 64GB memory with Ubuntu 18.04 LTS. The network bandwidth between machines is up to 10 GB/s. Among the 11 machines, one is as the user and the honest broker for SMCQL-GIS and Conclave-GIS, and the other 10 are data silos.

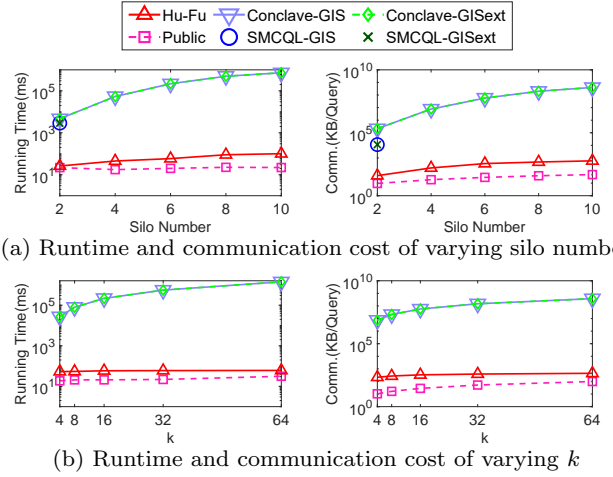


Fig. 6: Performance of federated kNN query.

7.2.1 Performance of Federated kNN Query

Fig. 6a shows the runtime and communication cost of (asymmetric) federated kNN query. Hu-Fu is 109.6 \times to 7,198.8 \times faster than SMCQL-GIS and Conclave-GIS, and has 2 to 5 orders of magnitude lower communication cost. When the number of silos increases from 2 to 10, the runtime and communication cost of Hu-Fu only increase by up to 2.9 \times and 13.9 \times , while those of Conclave-GIS drastically increase by up to 153.3 \times and 1,884.3 \times . Both metrics of Hu-Fu increase since the secure comparison and set union used in this query grow linearly with the silo number. Compare with Conclave-GIS and SMCQL-GIS, the runtime and communication cost of Conclave-GISext and SMCQL-GISext marginally increase (less than 20 ms and 200 KB respectively), which shows that our secure set union can efficiently assemble query results without an honest broker.

We also vary k from 4 to 64 and plot the running time and communication cost in Fig. 6b. As k increases from 4 to 64, the running time and communication cost of Hu-Fu only increase by 0.1 \times and 1.1 \times , while those of Conclave-GIS increase by 51.3 \times and 50.7 \times . The impact of k is less obvious than the silo number on Hu-Fu, because only the secure set union is linearly dependent on k . Again, the efficiency of Conclave-GISext is similar to that of Conclave-GIS. The drastic increase in running time and communication cost of Conclave-GIS and Conclave-GISext is expected because it involves many secure primitives that are time-consuming.

To show the improvement of DP optimization in kNN queries, we list the percentage of running time and communication cost reduced by DP in Table 4. With DP, the running time is reduced by up to 19.6%, and the communication cost by up to 47.7%. Compared

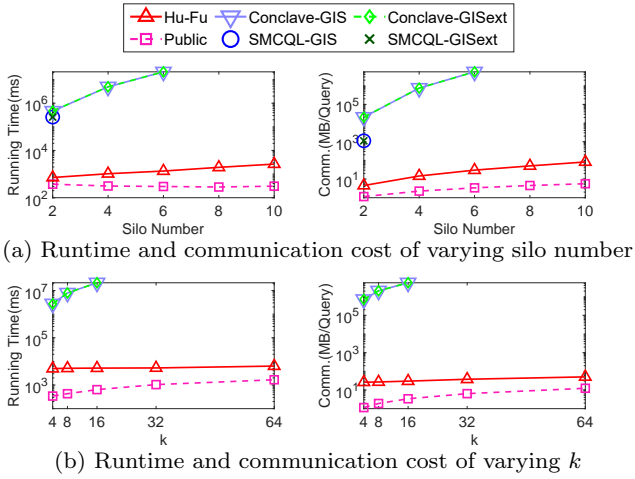


Fig. 7: Performance of federated kNN join.

with the improvement, the overhead of injecting the DP noise is very marginal, which takes $2 \mu s$ time cost and less than 1 KB communication cost when processing one federated kNN query. Such a notable improvement is because the complexity of DP noise injection is $O(1)$ and the summation only requires for transmission of n integers, while a secure comparison has $O(n)$ time complexity and $O(n^2)$ communication cost.

7.2.2 Performance of Federated kNN Join

Fig. 7a shows the results of (asymmetric) federated kNN join. Results of Conclave-GIS and Conclave-GISext with 8-10 silos are omitted since they incur over 6 hours for a single query. Hu-Fu is the most efficient, which is up to $360.2\times$ and $15,814.2\times$ faster than SMCQL-GIS and Conclave-GIS with $247.8\times$ and $185,151.0\times$ lower communication cost. The time and communication cost of SMCQL-GISext and Conclave-GISext slightly increase over SMCQL-GIS and Conclave-GIS.

Fig. 7b illustrates the impact of k . As k increases above 32, Conclave-GIS and Conclave-GISext require longer than 6 hours to process a federated kNN query. Thus, we can only provide their partial results (when $k \leq 16$). In contrast, Hu-Fu demonstrates superior advantages in the efficiency, achieving at least $553\times$ faster and $27,404\times$ lower communication cost compared to Conclave-GIS. As k increases from 4 to 64, the running time and communication cost of Hu-Fu rise by 28% and 48% respectively. The experimental trends for the federated kNN join closely resemble those observed in the federated kNN query, since a federated kNN join is decomposed into multiple federated kNN queries for all the compared solutions.

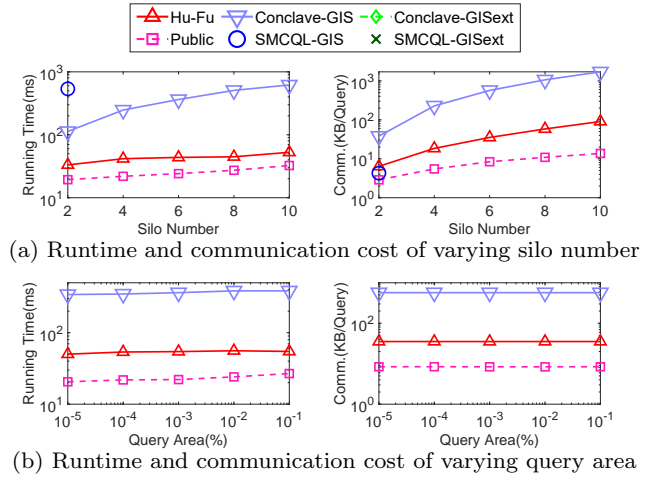


Fig. 8: Performance of federated range counting.

7.2.3 Performance of Federated Range Counting

Fig. 8 shows the results of (asymmetric) federated range counting. This query only returns the counting result and thus does not need a secure set union to protect data ownership. Hence, we exclude SMCQL-GISext and Conclave-GISext since they only differ from SMCQL-GIS and Conclave-GIS with an extra secure set union, which is unnecessary in this query. Hu-Fu is up to $15.2\times$ faster than SMCQL-GIS with a slightly higher communication cost (within 7 KB). Considering the increasing network bandwidth, the gap in communication cost is acceptable. Compared with Conclave-GIS, Hu-Fu is up to $10.8\times$ faster with $17.9\times$ lower communication cost. The running time and communication cost of Hu-Fu increase by $0.6\times$ and $13.2\times$ respectively when silo number increases to 10, mainly due to the secure summation.

We also demonstrate the impact of the query area on query efficiency in Fig. 8b. As is shown, the running time of all methods is relatively stable. It is expected because secure operations are the bottleneck of running time whereas the larger query area only increases the running time of plaintext operations.

7.2.4 Performance of Federated Range Query

Fig. 9 illustrates the results of (asymmetric) federated range query. The efficiency of SMCQL-GIS and Conclave-GIS is the same as Public (*i.e.*, the non-secure baseline), because they both rely on an honest broker to securely collect partial answers in each silo without leaking them to any others. Under this assumption, all systems can be reduced to Public, which uses a server (*e.g.*, an honest broker in SMCQL-GIS and a center server in Public) to directly collect local range query result from each

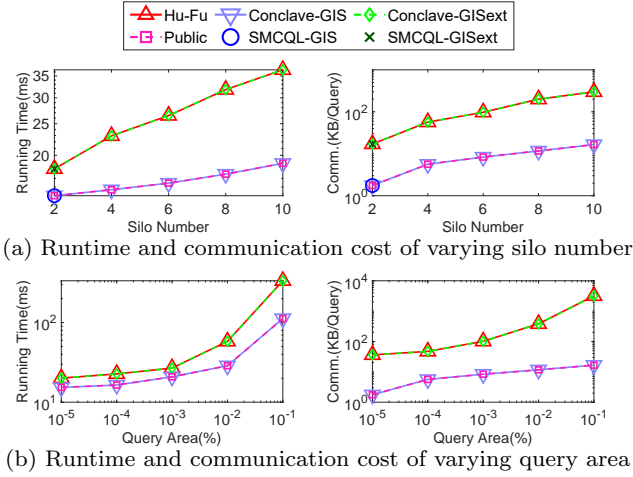


Fig. 9: Performance of federated range query.

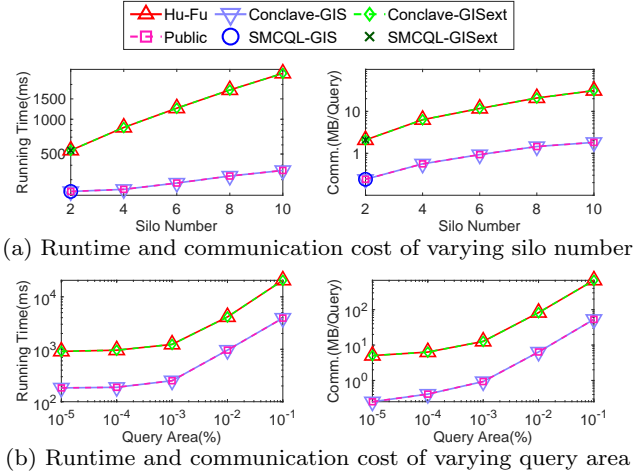


Fig. 10: Performance of federated distance join.

silos. For example, Hu-Fu with an honest broker also has the same efficiency as Public (see Appendix E).

Under a more general setting without this assumption, Hu-Fu, SMCQL-GISext and Conclave-GISext have the same efficiency because they all use our secure set union for results assembling. The usage of secure set union only leads to a marginal increase in running time (within 250 ms) and communication cost (lower than 3.1 MB) over Public. Note that the order of increase in running time and communication cost matches the complexity analysis for the secure set union in Sec. 5.2, which grows linearly with the silo number and the size of data returned. As shown in Fig. 9b, when the query area expands, all methods have a higher running time and communication cost, due to the increase of the number of spatial objects in the final result.

7.2.5 Performance of Federated Distance Join

Fig. 10 presents the performance of (asymmetric) federated distance join. Note that all the methods treat federated distance join as multiple independent federated range queries, where the total number of these range queries is $|R| = 100$ in this test. Thus, it is reasonable that the ranking of all the methods is similar to that in federated range query (see Fig. 9). The result of federated distance join when varying the query area shows a similar pattern with that of federated range query. This is because a federated distance join is decomposed into multiple federated range queries for all the compared solutions. The increase of both running time and communication cost is caused by the increase of the number of retrieved spatial objects.

7.2.6 Summary of Major Findings

We have observed the following findings in the experiments of asymmetric queries.

- Hu-Fu is up to $15,814.2\times$ faster than SMCQL-GIS and Conclave-GIS, with up to 5 orders of magnitude lower communication cost. The efficiency gain of Hu-Fu over the baselines is more notable in federated kNN query, kNN join, and range counting, which is at least $2.4\times$ faster in time cost and $4.9\times$ lower in communication cost than Conclave-GIS.
- SMCQL-GIS and Conclave-GIS are more efficient in federated range query and distance join, because these baselines are reduced to Public and need no secure operation with the honest broker. Note that for federated range query and distance join, Hu-Fu achieves the same efficiency as SMCQL-GISext and Conclave-GISext, the variants of SMCQL-GIS and Conclave-GIS without an honest broker.
- The experimental trends of federated kNN join and distance join are similar to those of federated kNN query and range query for all compared solutions. This is reasonable since a federated kNN join or distance join is decomposed into a series of federated kNN queries or range queries.

7.3 Experiments on Symmetric Queries

Parameter Setting. The parameter configurations for the query workloads are identical to those introduced in Section 7.2. Beyond these parameters, the privacy budget ϵ also affects the query performance of Hu-Fu. In general, a smaller ϵ (*i.e.*, stricter privacy preservation) leads to lower efficiency than a larger ϵ . Please refer to Appendix I for the evaluation of varying ϵ in

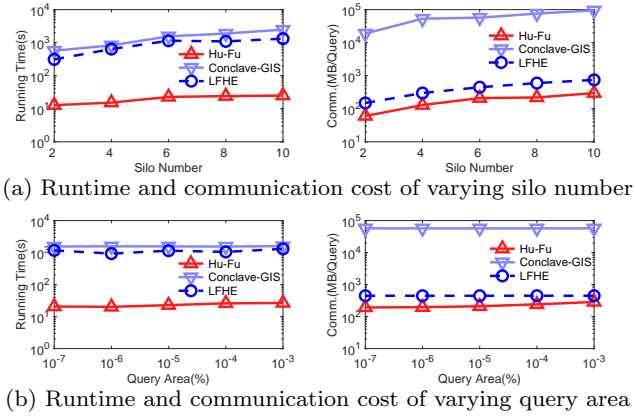


Fig. 11: Performance of federated range counting.

Hu-Fu. We have omitted reporting the results of the federated kNN join, since the query efficiency of Hu-Fu’s federated kNN join can be inferred from the results of its federated kNN query, as evidenced by previous experiments. We also omit the results of SMCQL-GIS, SMCQL-GISext, and Conclave-GISext, since their results are similar to Conclave-GIS when answering symmetric queries. Additional baselines, LFHE [37] and PINED-RQ++ [48], are involved in this experiment, where PINED-RQ++ is limited to federated range query and distance join.

Environment. Due to the expired funding support, the hardware environment for testing asymmetric queries is no longer accessible, so all solutions to symmetric queries are tested in a new hardware environment. Specifically, this new environment is composed of 5 cloud servers. Each server has Intel Xeon(R) Platinum 8361HC CPU 2.60GHZ processors and 32GB memory with Ubuntu 18.04 LTS OS. The network bandwidth between servers is up to 1.5 GB/s and may fluctuate slightly at different times. Four of the five servers act as data silos. Since this experiment requires up to 10 data silos, each of these four servers can host up to 3 data silos using different processes. The remaining server is only used as the user’s client, facilitating parallel data transmission between data silos for compare solutions. For each query type, we generate 50 queries, repeat 10 times for each query, and report the average results.

7.3.1 Performance of Federated Range Counting

Fig. 11 shows the results of (symmetric) federated range counting. Hu-Fu always outperforms the compared baselines in both running time and communication cost. Hu-Fu takes up to $98.9\times$ shorter runtime and up to $410.7\times$ lower communication cost than Conclave-GIS.

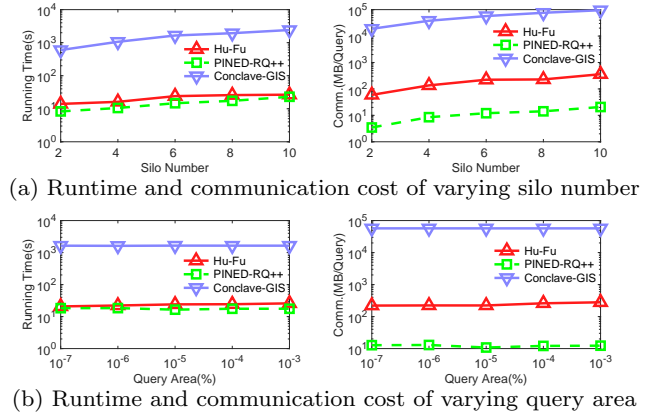


Fig. 12: Performance of federated range query.

Hu-Fu is also up to $56.6\times$ faster with up to $2.7\times$ lower communication cost compared to LFHE.

We can also observe that as the silo number increases from 2 to 10, the running time and communication cost of all compared algorithms tend to increase. This pattern is reasonable, as more data silos require more secure computations among them. When the query area expands from small to large, the runtime and communication overhead of Hu-Fu increase slightly. This trend in Hu-Fu is attributable to the increase of candidates for secure distance comparisons.

7.3.2 Performance of Federated Range Query

Fig. 12 illustrates the performance of (symmetric) federated range queries. The results of LFHE are not reported because the running time is over 1 hour and the memory consumption exceeds the server configuration (32 GB). In fact, LFHE can only handle a small-scale dataset. For example, LFHE already takes 13 minutes and 3 MB of communication cost when the data size is 10^4 . Although PINED-RQ++ is more efficient than Hu-Fu and Conclave-GIS, it suffers from *two major drawbacks*: (1) it leaks locations of the candidate data objects that are not part of the true answer to the query user and (2) it is only capable of retrieving approximation results. For example, the recall of PINED-RQ++ is 72.8%-96.8% when varying the silo number and 72.2%-92.3% when varying the query area.

Aside from PINED-RQ++, Hu-Fu achieves the best performance in the efficiency and security. It takes at least $42.7\times$ shorter time with at least $204.2\times$ lower communication cost compared to Conclave-GIS. Fig. 12 exhibits a similar trend in the efficiency variation of Hu-Fu when comparing to the results in Fig. 11 (for federated range counting). This similarity arises because the query decomposition plans for symmetric fed-

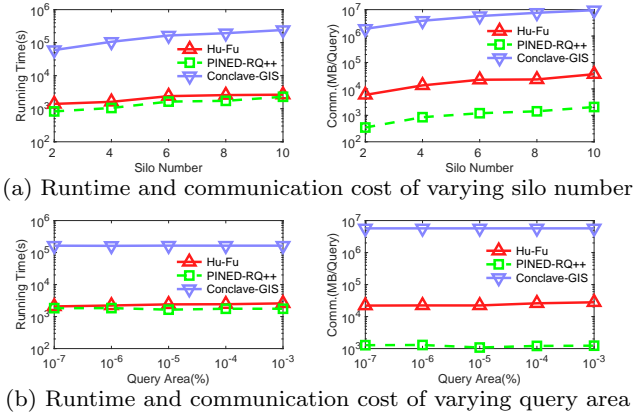


Fig. 13: Performance of federated distance join.

erated range queries and counting are quite identical. The main difference lies in the additional secure set union operator required for federated range queries.

7.3.3 Performance of Federated Distance Join

Fig. 13 presents the performance of (symmetric) federated distance joins. Here, a symmetric federated distance join employs the same strategy as the asymmetric federated distance join, treating the join operation as multiple independent federated range queries, with $|R| = 100$. Therefore, the relative performance of all methods here is similar to that of symmetric federated range query (see Fig. 12). Compared with PINED-RQ++, although Hu-Fu incurs higher communication cost, it takes nearly the same amount of running time and offers a more secure solution that produces the exact query answer. Compared to Conclave-GIS, Hu-Fu is up to 90.7 \times faster and up to 332.5 \times lower in communication cost. By contrast, LFHE is the least efficient method and cannot respond to a join query within 6 hours, so we cannot report its result in Fig. 13.

7.3.4 Performance of Federated kNN Query

Fig. 14 shows the running time and communication cost of (symmetric) federated kNN queries. We cannot report some results of LFHE and Conclave-GIS because they cannot terminate within 1 hour to process a single query or the maximum memory usage is beyond the limitation (32 GB) of the cloud server. Additionally, PINED-RQ++ is unable to handle KNN queries, so it is excluded in this evaluation. According to the experimental results, Hu-Fu is the most efficient solution to federated kNN queries. For instance, Hu-Fu is up to 551.8 \times faster than Conclave-GIS and 56.2 \times faster than

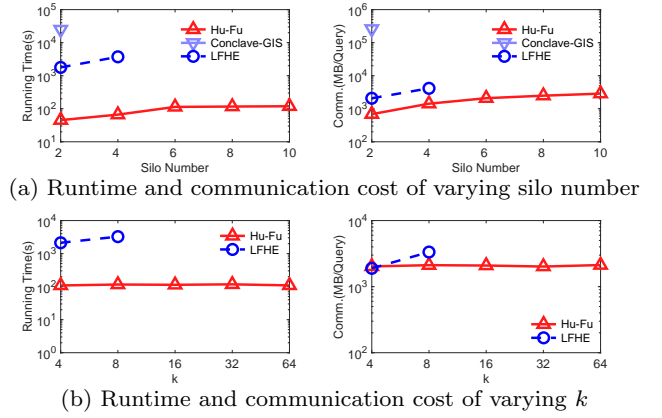


Fig. 14: Performance of federated kNN query.

Table 5: Improvement with upper bound in Lemma 1

	k	4	8	16	32	64
Running time		6.3 \times	4.2 \times	4.0 \times	3.6 \times	3.4 \times
Communication		5.0 \times	3.9 \times	3.6 \times	3.4 \times	3.1 \times

LFHE, with up to 380.0 \times and 3.0 \times lower communication overhead than them, respectively. Similar to previous results, the time cost of Hu-Fu gradually increases as the silo number increases. By contrast, the efficiency of Hu-Fu does not notably change as k increases.

We also evaluate the improvement of using the upper bound in Lemma 1. As shown in Table 5, this optimization improves the running time by up to 6.3 \times and reduces the communication cost by up to 5.0 \times .

7.3.5 Summary of Major Findings

The major findings in the experiments of symmetric queries are summarized as follows.

- Hu-Fu is at least 42.7 \times faster than Conclave-GIS, with at least 204.2 \times lower communication overhead. Compared to LFHE, Hu-Fu is up to 56.6 \times faster with up to 3.0 \times lower communication overhead.
- Although PINED-RQ++ is more efficient than the others, it suffers from three significant drawbacks: (1) violations on the data privacy, (2) inability to provide exact answers, and (3) limited support to only federated range queries. By contrast, our Hu-Fu can address all these drawbacks effectively.
- When comparing with the evaluations of asymmetric queries in Sec. 7.2, it is evident that running time and communication overhead both escalate when Hu-Fu or Conclave-GIS processes symmetric queries. Consequently, symmetric queries pose a greater challenge than asymmetric queries, primarily due to the additional concern for query privacy.

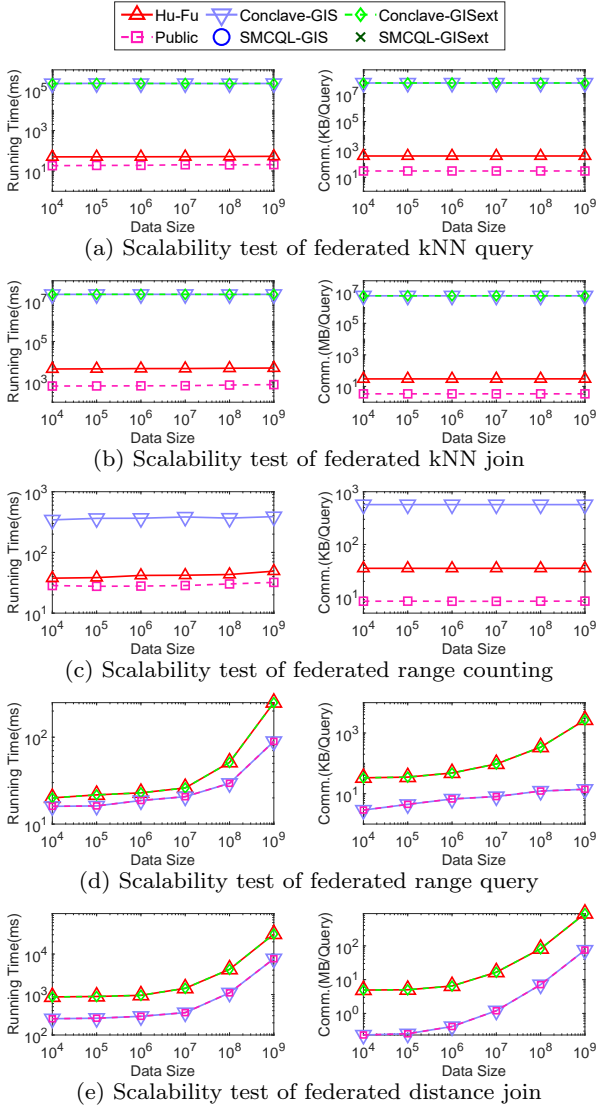


Fig. 15: Scalability test on asymmetric queries.

Remark. To assess the impact of query privacy on time efficiency, we can compare the running time of asymmetric and symmetric queries in the new hardware environment. Our evaluation shows that the baseline Conclave-GIS takes $4,054\times$ longer to protect query privacy in federated range queries, while our solution Hu-Fu reduces this gap to $59\times$. Please see Appendix H for more detailed results.

7.4 Experiments on Scalability Tests

In the following, we report the results of scalability tests on asymmetric queries and symmetric queries in Sec. 7.4.1 and Sec. 7.4.2, respectively.

7.4.1 Scalability Test on Asymmetric Queries

Parameter Setting. In the following, we scale the total number of spatial objects from 10^4 to 10^9 over OSM dataset to assess the scalability of Hu-Fu. Other parameters are set to the default values as in Sec. 7.2. For example, the number of silos is 6, $k = 16$ for federated kNN query and kNN join, and the query area for federated range query, range counting and distance join is 0.001%. Since SMCQL-GIS and SMCQL-GISext only support two silos, they are excluded in the scalability test. The running time and communication cost on all five spatial queries are shown in Fig. 15.

Result and Analysis. For a fixed data size, we observe that Hu-Fu is notably more efficient than Conclave-GIS and Conclave-GISext on federated kNN query, kNN join and range counting (see Fig. 15a-15c). For federated range query and distance join, Conclave-GIS behaves the same as Public due to the honest broker, while Hu-Fu achieves the same efficiency as Conclave-GISext, which requires no honest broker.

We are more interested in the efficiency with the increase of data size. We observe that the efficiency of federated kNN query, kNN join and range counting is insensitive to the increase of the data size. This is because the increase of data size mainly affects the time cost of plaintext primitives, which only accounts for a small portion (due to efficient indexes in each silo) in the running time. In contrast, the running time and communication cost of federated range query and distance join notably increase with the increase of the data size because more spatial objects are retrieved in each silo, which leads to a higher cost for both plaintext range query and secure set union.

Takeaways. Hu-Fu trivially scales with data size for federated kNN query, kNN join and range counting, because these queries are relatively insensitive to data size. Both metrics of Hu-Fu increase with the data size for federated range query and distance join, yet Hu-Fu is still reasonably efficient for them on large-scale data. For example, in Hu-Fu, an asymmetric federated range query takes 250 ms running time and 2.6 MB communication cost on the data size of 10^9 .

7.4.2 Scalability Test on Symmetric Queries

Parameter Setting. As for symmetric queries, we evaluate the scalability of Hu-Fu by scaling the OSM dataset, gradually increasing the total number of spatial objects from 10^4 to 10^8 . The baseline selection here is identical to that in Sec. 7.3. Although we omit the results of spatial joins, the scalability of the federated distance join and kNN join can be inferred from the scalability of the

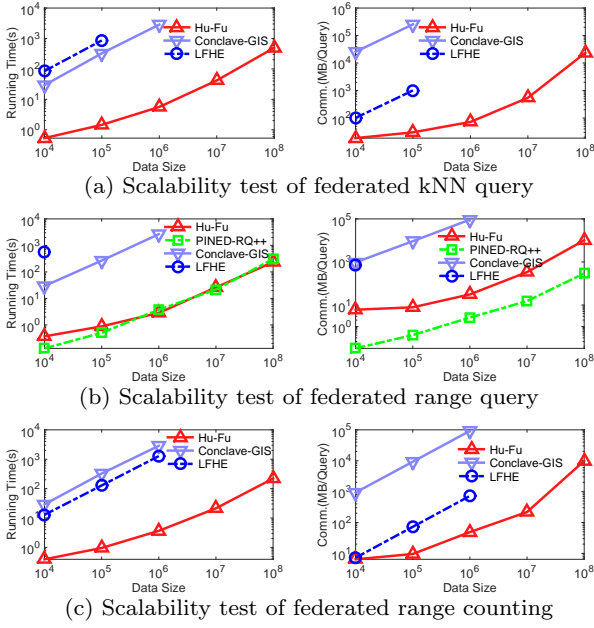


Fig. 16: Scalability test on symmetric queries.

federated range query and kNN query, as demonstrated in the previous experiments in Sec. 7.2 and Sec. 7.3.

Result and Analysis. As shown in Fig. 16, for any data size, Hu-Fu significantly outperforms Conclave-GIS and LFHE in both running time and communication overhead for all the tested queries. For example, in Fig. 16a, Hu-Fu is at least $54.8\times$ and $1375.7\times$ faster than Conclave-GIS and LFHE, respectively. Neither Conclave-GIS nor LFHE can efficiently process these queries over large-scale datasets (*e.g.*, when the data size is over 10^6). After running for more than 6 hours, neither of them have terminated, so we are unable to obtain their full results. Moreover, when processing symmetric federated range queries, both Conclave-GIS and LFHE take at least 2 orders of magnitude higher communication cost than Hu-Fu.

By contrast, PINED-RQ++ is slightly faster than Hu-Fu (no more than $2.8\times$) and has lower communication overhead in Fig. 16b. However, since PINED-RQ++ may violate the data privacy during the query processing and can only obtain approximate results, the performance gap is acceptable. For instance, the recall of PINED-RQ++ can be lower than 80%, and such a low recall may lead to an unsatisfying service experience in our motivation scenarios like contact tracing, where accurate results are crucial.

Takeaways. Unlike the scalability tests for asymmetric queries, the efficiency of symmetric federated range query, range counting, and kNN query is highly sensitive to the data size. Specifically, when processing *asymmetric queries*, all compared solutions leverage plain-

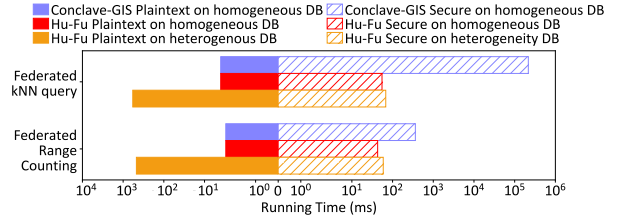


Fig. 17: Running time breakdown.

text operators to filter most of the data objects. For example, a plaintext range query can ensure no false positive candidate for asymmetric federated range queries, while a plaintext kNN query results in only k false positive candidates for asymmetric federated kNN queries. Consequently, the number of false positive candidates is (almost) independent of data size.

However, when processing *symmetric queries* that require additional protections for the query location, although plaintext operators are still used to reduce the candidate size, no solution can ensure a constant number of false positive candidates. Even in our Hu-Fu, the number of candidate objects awaiting secure verification after filtering by plaintext operators is proportional to the data size. This contributes to the increased time and communication overhead required to process symmetric queries compared to asymmetric ones.

7.5 Experiments on Heterogeneous Data Silos

This experiment aims to demonstrate the feasibility of Hu-Fu on heterogeneous spatial databases. Specifically, we use 6 different databases for each silo on the BJ dataset: PostGIS [10], MySQL [8], SpatiaLite [11], Simba [64], GeoMesa [27], and SpatialHadoop [21]. Other parameters are set as the default values as in Sec. 7.2.

Fig. 17 plots the running time breakdown *i.e.*, plaintext vs. secure primitives for radius-unknown (*i.e.*, asymmetric federated kNN query) and radius-known (*i.e.*, asymmetric federated range counting) queries. We make the following observations.

- Given homogeneous underlying spatial databases (*i.e.*, PostGIS), our Hu-Fu significantly reduces the running time of secure primitives *e.g.*, $3,935.4\times$ compared with Conclave-GIS for federated kNN query. Such acceleration in secure primitives is the primary contributor to Hu-Fu’s gain in running time.
- Heterogeneous underlying spatial databases affect the running time. Specifically, the running time of plaintext primitives is limited by the slowest spatial database, which may increase the overall query processing time. In this experiment, the running time

of plaintext primitives notably increases from 4 ms to 579 ms when replacing PostGIS with heterogeneous databases (where SpatiaLite and MySQL are the slowest). It takes even longer time than the secure primitives in Hu-Fu. The running time of secure primitives also marginally increases, due to idle waiting for the local results from the slowest silo.

For other symmetric federated spatial queries (*i.e.*, range query, distance join, and kNN join), please refer to Appendix D for their detailed results.

Takeaways. Hu-Fu functions with data silos running heterogeneous databases. Although Hu-Fu dramatically speeds up the secure primitives in a federated spatial query, the efficiency of plaintext primitives in each silo’s databases may affect the overall running time. Particularly, the time cost of plaintext primitives can be limited by the slowest database in the federation. To unleash the full potential of Hu-Fu, fast spatial databases in each silo are recommended.

8 Related Work

Distributed spatial database systems are popular solutions to query processing on big spatial data. These systems improve query processing via data partition and indexing techniques (*e.g.*, R-tree [44]) in Hadoop (*e.g.*, SpatialHadoop [21]) or Spark (*e.g.*, Simba [64]). However, the data partition techniques are inapplicable in a data federation since the entire data is held by the autonomous data silos. Moreover, security is not the major concern in these systems.

Past studies of secure spatial query processing mainly focus on encrypted databases [32], where data is encrypted and stored in a third-party platform (*e.g.*, a cloud platform) to process queries securely. For example, existing work [22, 37, 62, 65] study the secure kNN query on encrypted databases and prior studies [48, 53, 58, 63] focus on securely processing range queries. In these studies, a data owner outsources its data and hence the sensitive data is encrypted before being uploaded to a third party. Intuitively, homomorphic encryption techniques (*e.g.*, Paillier and SEAL [12]) are used to guarantee security. Different from this setting, in a data federation, data silos autonomously manage their own data and hence do not need to encrypt their own data and upload it to a third party. Besides, our experiments demonstrate that extending these solutions [37, 48] to the scenario of the data federation can be either insecure or inefficient.

Rather than the general distributed databases or outsourced databases, our work is more aligned with the problem settings of federated databases and data

federation, where the entire dataset is held in multiple autonomous databases. Early research on federated databases focused on finding solutions to access data in autonomous databases [45], while recent studies on federated databases support diverse data types, *e.g.*, on federated graph databases [57]. Note that the autonomous database here means that data can be only managed by its held silo.

Data federation is an emerging concept developed from federated databases. It shares a similar architecture with federated databases. Yet, the *major difference* is that a data federation imposes certain secure requirements during query processing, while a federated database does not. For example, SMCQL [14] is probably the first secure query processing solution over a data federation and Conclave [56] is the state-of-the-art solution. More recent studies explored efficient solutions to join queries [33, 40, 59] in a data federation. All these studies adopt SMC techniques to achieve secure query processing for *relational* data with *exact* results.

Exact federated queries have been explored for various data types. Our preliminary work [51] and its accompanying demonstration system [46] focus on exact federated queries over spatial data federation. Zhang *et al.* [70] propose an efficient method that leverages the Intel SGX to securely perform similarity searches over a data federation under metric spaces. For example, the metric distance can be the graph edit distance for graph data or the edit distance for sequence data.

Existing studies also investigate *approximate* query processing over a data federation. SAQE [16], Cryptic [20], and Shrinkwrap [15] use differential privacy to trade off between accuracy and efficiency in processing relational queries. Others study approximate kNN queries [69] and range counting [38, 49] over a spatial data federation. In contrast, we focus on *exact* query processing, since accurate results can be crucial for spatial applications like contact tracing [26].

In short, our work is inspired by the emerging trend of secure query processing over a data federation, yet focuses on spatial queries with exact results. Our Hu-Fu significantly improves the efficiency of federated spatial queries over the extensions of SMCQL [14] and Conclave [56], the state-of-the-arts for relational data. Moreover, unlike most existing studies that solely focus on protecting data privacy under this emerging scenario, Hu-Fu also considers the application need for preserving the query privacy. For example, in applications like contact tracing, spatial queries often contain sensitive location data of patients, thereby necessitating the protection on the query privacy.

9 Conclusion

In this paper, we propose the first system Hu-Fu for efficient and secure spatial queries over a data federation. Existing solutions are inefficient to process such queries due to excessive secure distance operations and the usage of general-purpose secure multi-party computation (SMC) libraries for implementing secure operators. To overcome the inefficiency, we design a novel query rewriter to decompose the spatial queries into as many plaintext operators and as few secure operators as possible. In particular, our secure operators have dedicated implementations faster than general-purpose SMC libraries. Moreover, Hu-Fu supports heterogeneous spatial databases (*e.g.*, PostGIS, Simba, GeoMesa, and SpatialHadoop), as well as query input in native SQL. Finally, extensive experiments show that Hu-Fu is up to 4 orders of magnitude faster and takes 5 orders of magnitude lower communication cost than the state-of-the-arts. In the future study, we plan to support more spatial queries, *e.g.*, spatial keyword search.

References

- (2022) Facebook scandal: Who is selling your personal data? URL <https://www.bbc.com/news/technology-44793247>
- (2024) Acxiom. URL <https://www.acxiom.com>
- (2024) AMAP. URL <https://mobility.amap.com>
- (2024) China Mobile. URL <https://www.chinamobileltd.com>
- (2024) China Telecom. URL <http://www.chinatelecom-h.com>
- (2024) Communication travel card. URL <https://xc.caict.ac.cn/>
- (2024) Hu-Fu: Efficient and secure spatial queries over data federation. URL <https://github.com/BUAA-BDA/OpenHuFu/>
- (2024) MySQL. URL <https://www.mysql.com>
- (2024) OpenStreetMap. URL <https://www.openstreetmap.org>
- (2024) PostGIS. URL <https://www.postgis.org>
- (2024) SpatiaLite. URL <https://www.gaia-gis.it/fossil/libspatialite/index>
- Acar A, Aksu H, Uluagac AS, Conti M (2018) A survey on homomorphic encryption schemes: Theory and implementation. *ACM Comput Surv* 51(4):79:1–79:35
- Andrés ME, Bordenabe NE, Chatzikokolakis K, Palamidessi C (2013) Geo-indistinguishability: differential privacy for location-based systems. In: *CCS*, pp 901–914
- Bater J, Elliott G, Eggen C, Goel S, Kho AN, Rogers J (2017) SMCQL: secure query processing for private data networks. *PVLDB* 10(6):673–684
- Bater J, He X, Ehrich W, Machanavajjhala A, Rogers J (2018) ShrinkWrap: Efficient SQL query processing in differentially private data federations. *PVLDB* 12(3):307–320
- Bater J, Park Y, He X, Wang X, Rogers J (2020) SAQE: practical privacy-preserving approximate query processing for data federations. *PVLDB* 13(11):2691–2705
- Bettini C, Jajodia S, Samarati P, Wang XS (eds) (2009) *Privacy in Location-Based Applications, Research Issues and Emerging Trends*. Springer
- Bogdanov D, Laur S, Willemson J (2008) Sharemind: A framework for fast privacy-preserving computations. In: *ESORICS*, pp 192–206
- Calcite (2024) URL <https://calcite.apache.org/>
- Chowdhury AR, Wang C, He X, Machanavajjhala A, Jha S (2020) Crypte: Crypto-assisted differential privacy on untrusted servers. In: *SIGMOD*, pp 603–619
- Eldawy A, Mokbel MF (2015) SpatialHadoop: A mapreduce framework for spatial data. In: *ICDE*, pp 1352–1363
- Elmehdwi Y, Samanthula BK, Jiang W (2014) Secure k-nearest neighbor query over encrypted data in outsourced environments. In: *ICDE*, pp 664–675
- Emekçi F, Sahin OD, Agrawal D, Abbadi AE (2007) Privacy preserving decision tree learning over multiple parties. *Data Knowl Eng* 63(2):348–361
- Evans D, Kolesnikov V, Rosulek M (2018) A pragmatic introduction to secure multi-party computation. *Foundations and Trends in Privacy and Security* 2(2-3):70–246
- Facebook (2024) URL <https://www.facebook.com>
- Ferretti L, Wymant C, Kendall M, Zhao L, Nurtay A, Abeler-Dörner L, Parker M, Bonsall D, Fraser C (2020) Quantifying sars-cov-2 transmission suggests epidemic control with digital contact tracing. *Science* 368(6491):eabb6936
- GeoMesa (2024) URL <https://www.geomesa.org/>
- Gertz M, Jajodia S (eds) (2008) *Handbook of Database Security - Applications and Trends*. Springer
- Gkoulalas-Divanis A, Bettini C (eds) (2018) *Handbook of Mobile Data Privacy*. Springer
- Goldreich O (2009) *Foundations of cryptography: volume 2, basic applications*. Cambridge university press
- Google Maps (2024) URL <https://www.google.cn/maps/>

32. Hacigümüs H, Iyer BR, Li C, Mehrotra S (2002) Executing SQL over encrypted data in the database-service-provider model. In: SIGMOD, pp 216–227
33. Han F, Zhang L, Feng H, Liu W, Li X (2022) Scape: Scalable collaborative analytics system on private database with malicious security. In: ICDE, pp 1740–1753
34. Instagram (2024) URL <https://www.instagram.com>
35. Jurczyk P, Xiong L (2011) Information sharing across private databases: Secure union revisited. In: PASSAT/SocialCom, pp 996–1003
36. Keller M (2020) MP-SPDZ: A versatile framework for multi-party computation. In: CCS, pp 1575–1590
37. Kesarwani M, Kaul A, Naldurg P, Patranabis S, Singh G, Mehta S, Mukhopadhyay D (2018) Efficient secure k-nearest neighbours over encrypted data. In: EDBT, pp 564–575
38. Li M, Zeng Y, Chen L (2023) Efficient and accurate range counting on privacy-preserving spatial data federation. In: DASFAA, pp 317–333
39. Li N, Lyu M, Su D, Yang W (2016) Differential Privacy: From Theory to Practice. Morgan & Claypool Publishers
40. Li S, Zeng Y, Wang Y, Zhong Y, Zhou Z, Tong Y (2024) An experimental study on federated equi-joins. IEEE Trans Knowl Data Eng 36(9):4443–4457
41. Lindell Y (2021) Secure multiparty computation. Communications of the ACM 64(1):86–96
42. Liu C, Wang XS, Nayak K, Huang Y, Shi E (2015) ObliVM: A programming framework for secure computation. In: S&P, pp 359–376
43. Liu F, Zheng Z, Shi Y, Tong Y, Zhang Y (2024) A survey on federated learning: a perspective from multi-party computation. Frontiers Comput Sci 18(3):181,336
44. Mamoulis N (2011) Spatial Data Management. Synthesis Lectures on Data Management, Morgan & Claypool Publishers
45. Özsu MT, Valduriez P (2020) Principles of Distributed Database Systems, 4th Edition. Springer
46. Pan X, Tong Y, Xue C, Zhou Z, Du J, Zeng Y, Shi Y, Zhang X, Chen L, Xu Y, Xu K, Lv W (2022) Hu-fu: A data federation system for secure spatial queries. PVLDB 15(12):3582–3585
47. Parliament E, of the European Union TC (2024) The general data protection regulation (gdpr). URL <https://eugdpr.org>
48. Sahin C, Allard T, Akbarinia R, Abbadi AE, Pacitti E (2018) A differentially private index for range query processing in clouds. In: ICDE, pp 857–868
49. Shi Y, Tong Y, Zeng Y, Zhou Z, Ding B, Chen L (2023) Efficient approximate range aggregation over large-scale spatial data federation. IEEE Trans Knowl Data Eng 35(1):418–430
50. Sun N, Wang W, Tong Y, Liu K (2024) Blockchain based federated learning for intrusion detection for internet of things. Frontiers Comput Sci 18(5):185,328
51. Tong Y, Pan X, Zeng Y, Shi Y, Xue C, Zhou Z, Zhang X, Chen L, Xu Y, Xu K, Lv W (2022) Hu-Fu: Efficient and secure spatial queries over data federation. PVLDB 15(6):1159–1172
52. Tong Y, Zeng Y, Zhou Z, Liu B, Shi Y, Li S, Xu K, Lv W (2023) Federated computing: Query, learning, and beyond. IEEE Data Eng Bull 46(1):9–26
53. Tran HV, Allard T, d’Orazio L, Abbadi AE (2019) Range query processing for monitoring applications over untrustworthy clouds. In: EDBT, pp 666–669
54. TripAdvisor (2024) URL <https://www.tripadvisor.com>
55. Voigt P, Von dem Bussche A (2017) The EU General Data Protection Regulation (GDPR): A Practical Guide, vol 10. Springer
56. Volgushev N, Schwarzkopf M, Getchell B, Varia M, Lapets A, Bestavros A (2019) Conclave: secure multi-party computation on big data. In: EuroSys, pp 3:1–3:18
57. Vu X, Ait-Mlouk A, Elmroth E, Jiang L (2019) Graph-based interactive data federation system for heterogeneous data retrieval and analytics. In: WWW, pp 3595–3599
58. Wang P, Ravishankar CV (2013) Secure and efficient range queries on outsourced databases using rp-trees. In: ICDE, pp 314–325
59. Wang Y, Yi K (2021) Secure Yannakakis: Join-aggregate queries over private data. In: SIGMOD, pp 1969–1981
60. Wang Y, Zeng Y, Xu Y, Zhou Z, Tong Y (2024) Efficient and private federated trajectory matching. CoRR abs/2312.12012
61. Waze (2024) URL <https://www.waze.com>
62. Wong WK, Cheung DW, Kao B, Mamoulis N (2009) Secure knn computation on encrypted databases. In: SIGMOD, pp 139–152
63. Wu S, Li Q, Li G, Yuan D, Yuan X, Wang C (2019) ServeDB: Secure, verifiable, and efficient range queries on outsourced database. In: ICDE, pp 626–637
64. Xie D, Li F, Yao B, Li G, Zhou L, Guo M (2016) Simba: Efficient in-memory spatial analytics. In: SIGMOD, pp 1071–1085

65. Yao B, Li F, Xiao X (2013) Secure nearest neighbor revisited. In: ICDE, pp 733–744
66. Ye J (2019) Transportation: A data driven approach. In: KDD, p 3183
67. Yelp (2024) URL <https://www.yelp.com>
68. Zeng Y, Tong Y, Chen L (2019) Last-mile delivery made practical: An efficient route planning framework with theoretical guarantees. PVLDB 13(3):320–333
69. Zhang K, Tong Y, Shi Y, Zeng Y, Xu Y, Chen L, Zhou Z, Xu K, Lv W, Zheng Z (2023) Approximate k-nearest neighbor query over spatial data federation. In: DASFAA, pp 351–368
70. Zhang X, Wang Q, Xu C, Peng Y, Xu J (2024) Fed-knn: Secure federated k-nearest neighbor search. SIGMOD 2(1):V2mod011:1–V2mod011:26

A Application Scenarios of Query Privacy

Due to legal regulations (*e.g.*, GDPR [47]), protecting *data privacy* is now common in real-life scenarios, especially when spatial location implies the travel patterns or trajectories of a user. *Query privacy* is equally important has numerous application scenarios. Here are some detailed examples:

1. **Navigation:** Apps like Google Maps [31] and Waze [61] frequently query user locations to provide directions and traffic updates. Revealing exact locations can expose sensitive information about users' daily routines and personal lives.
2. **Social Networking:** Platforms like Facebook [25] and Instagram [34] allow users to check in at various locations. Protecting the privacy of these check-ins is crucial to prevent stalking and harassment.
3. **Location-based Advertising:** This service targets users based on their current or recent locations. Protecting user privacy ensures that advertisers do not have access to sensitive location data, which could be misused for targeted marketing.
4. **POI Search:** Applications like Yelp [67] and TripAdvisor [54] allow users to search for nearby points of interest like restaurants, cafes, or hotels. Ensuring the privacy of sensitive location data will prevent potential tracking and profiling of users' movement patterns by the platform.

Please refer to [17, 29] for more application scenarios of query privacy.

B Security Analysis

The security analysis of our query rewriter is based on the composition lemma in [30] (see its Sec. 7.3.1). The lemma states that given a secure protocol $\pi^{g|f}$ that can securely compute g based on a plaintext query f and a secure protocol π^f for operation f , the operation g can be securely computed by executing protocol $\pi^{g|f}$ but substitutes every plaintext query f with an execution of π^f .

We prove the security of query rewriter for asymmetric federated queries (Appendix B.1) and symmetric federated queries (Appendix B.2) separately. Finally, we present a case study on attacking asymmetric federated range counting in Appendix B.3.

B.1 Security Analysis of Asymmetric Federated Queries

In the following, we conduct a thorough analysis of the security of asymmetric federated queries in Hu-Fu, categorizing them into two distinct kinds: *radius-known queries* and *radius-unknown queries*.

- **Security of Radius-Known Queries.** Federated range query, range counting and distance join belong to *radius-known queries*. For *federated range query* and *distance join*, the secure protocol $\pi^{g|f}$ is the secure set union (Def. 5) based on the query results of one or multiple plaintext range queries over the federation. The protocol π^f corresponds to the plaintext range query (Def. 1). The protocol π^f is secure, since the plaintext range query only occurs in each silo, which will not leak any information to other silos. From the composition lemma, federated range

query and distance join are secure against semi-honest adversaries. Federated range counting is also secure based on a similar proof (*i.e.*, replacing $\pi^{g|f}$ with the secure summation in Def. 3 and π^f with the plaintext range counting in Def. 1).

- **Security of Radius-Unknown Queries.** The federated kNN query and kNN join belong to *radius-unknown queries*. For federated kNN query (Alg. 1), let protocol π^f be the secure count comparison between k and the local counts. The security of π^f can be proved in the same way as before. Denote protocol $\pi^{g|f}$ as Alg. 1 based on the plaintext count comparison query. We then show protocol $\pi^{g|f}$ is also secure. Assuming semi-honest adversaries, each silo can simulate its view of the execution of the protocol. Specifically, knowing the range $[l, u]$ used at the beginning of a round, each silo can compute r used in that round. If r is the same as the final radius, it concludes that the protocol must have ended in this round. Otherwise, it simply updates the range to that side of r which contains the final radius. Along with the knowledge of the initial range, this shows that each silo can simulate the execution of the protocol, which means the protocol $\pi^{g|f}$ is secure. Thus, our Alg. 1 is secure against semi-honest adversaries based on the composition lemma. The security of federated kNN join can be proved similarly.

B.2 Security Analysis of Symmetric Federated Queries

Similar to the previous security analysis, we analyze the security of symmetric federated queries in Hu-Fu from two categories: *radius-known queries* and *radius-unknown queries*.

- **Security of Radius-Known Queries.** For *federated range query*, the secure protocol $\pi^{g|f}$ employs the secure distance comparison (Def. 6), which relies on the fully homomorphic encryption (FHE) scheme (*i.e.*, the BGV scheme [12]), and the secure set union (Def. 5). The protocol π^f corresponds to the plaintext range query (Def. 1). It is secure because the plaintext range query is executed within each silo individually, thereby preventing any information leakage to other silos. According to the composition lemma, the query rewriter for symmetric federated range queries ensures data privacy. Regarding the additional privacy requirement of query privacy, the secure protocol, secure location perturbation (Def. 7), is based on the BPL mechanism in [60]. This mechanism has been proven to satisfy (ϵ, δ) -Geo-Indistinguishability (Geo-I) [13], a widely-adopted privacy notion for preserving location privacy. Thus, the query privacy is due to the protection of query location with secure location perturbation. As for *federated distance join*, the data privacy and query privacy can be guaranteed, since a federated distance join is decomposed into a series of federated range queries, and the security of each federated range query has been proven. The security of *federated range counting* shares similarities with that of federated range query, primarily due to their near-identical decomposition plans (see Table 3). However, a key difference lies in the nature of their query answers. While federated range query requires a secure set union to collect spatial objects falling within the specified range, federated range counting simply returns an integer count as its output. This distinction underscores the minimal exposure of information in federated range counting, further reinforcing its security.

Table 6: A case study on hacking (asymmetric) federated range counting in Hu-Fu.

Colluded silos	Attack success rate	Running time
1	0	No result for 2 months
2	0	No result for 2 months
3	0	No result for 2 months

- **Security of Radius-Unknown Queries.** To prove the security of radius-unknown queries, we only need to show the security of federated kNN query, since a federated kNN join is decomposed into multiple independent federated kNN queries by Hu-Fu. Similar to the decomposition plan for asymmetric federated kNN query, the query rewriter for symmetric federated kNN query also utilizes a binary-search procedure. The key difference is that the initial upper bound is determined by plaintext kNN queries and secure location perturbation. Thus, for this step, the secure protocol $\pi^{g|f}$ is due to the privacy guarantee of (ϵ, δ) -Geo-Indistinguishability (Geo-I) [13]. The security analysis for the other steps is similar to that for asymmetric radius-unknown queries.

B.3 Case Study on Semi-Honest Attack

Case Study. We have also conducted a case study to show how Hu-Fu defenses against the semi-honest adversaries when processing asymmetric federated range counting.

- **Experimental Setup.** The case study is based on the BJ dataset with 6 silos (*i.e.*, our default setting) and we vary the number m of colluded adversaries (*i.e.*, silos) from 1 to 3. Then, when processing asymmetric federated range counting, the attack will occur at the only secure operator (*i.e.*, secure summation [23]). To attack the secure summation, we use the method introduced in [23]. In general, a secure operator successfully defends against the attack, if adversaries cannot get any feasible result after a long time (*e.g.*, two months in our experiment). Furthermore, the other experiment parameters are the same as those in Sec. 7.2.
- **Experimental Result.** As shown in Table 6, even when 50% of the silos collude, they cannot infer any feasible result even after two months’ cracking. Note that it is commonly assumed in existing work that the portion of colluded adversaries is less than 50%. Based on the computational security [28, 30, 41], this result proves Hu-Fu is hard to attack.

C Handling Ties in Symmetric Federated kNN Query

We provide a solution for handling ties in symmetric federated kNN query and conduct an experimental comparison. **Main Idea.** To break our assumption and handle the special case of ties, we need to address two questions: (1) how to identify such cases and (2) how to retrieve exactly k nearest neighbors, using the following ideas:

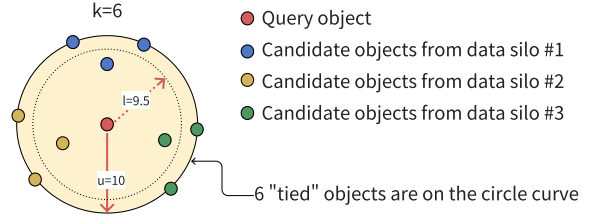


Fig. 18: Illustration of the special case of symmetric federated kNN query ($k = 6$).

- (1) **Identify such cases.** For any searching radius r during the binary-search, whenever the total number of spatial objects closer to the query object q than r is equal to k (*i.e.*, the result of secure count comparison is $sign = 0$), we can find exactly k nearest neighbors (*i.e.*, there are no other objects whose distances to q are also the k th nearest distances). When the binary-search terminates and $sign$ is always non-zero, we will encounter the mentioned case (we call it “ties” for brevity).
- (2) **Retrieve exactly k nearest neighbors.** Under this case, we will obtain the final lower bound l and upper bound u of the searching radius for k th nearest neighbor. Since the binary-search has been terminated, we can derive that $u - l$ must be smaller than the distance precision ϵ_0 . Now, we *first* use the circular range circle(q, l) to locate all spatial objects that are strictly inside the circle curve (*i.e.*, by a symmetric federated range query), and denote the total count as k_l . Notice that we can temporarily cache the partial answers in each data silo at this step. *Next*, we only need to randomly pick $k - k_l$ spatial objects from those objects on the circle curves. To achieve this goal, we sequentially request such spatial objects from every data silo, and set a limitation on how many spatial objects are indeed required. *Lastly*, a secure set union operator is used to collect the k nearest neighbors, since the partial answer has been cached in each data silo.

Alg. 4 illustrates the detailed procedure of using the main idea to process the cases of ties. Specifically, lines 3-9 represent the procedure of binary-search. In line 10, when $sign$ equals 0 in some iteration of the binary search, we know the current searching radius r is already enough to retrieve the exact k nearest neighbors. A symmetric federated range query can drive the query answer. By contrast, in line 12, $sign$ has never become 0, then we know we are facing the case of ties now. To handle this case, line 14 computes how many data objects are strictly inside the circle curve, denoted by k_l . Thus, we only need to pick $k - k_l$ from those tied objects on the circle curve. To do so, lines 16-20 sequentially request certain objects from each data silo, and limit the number of retrieved object no larger than $remain$ (since we do not need more). To avoid data objects are duplicate statistics, lines 14 and 17 will cache the (partial) answer locally at each silo. Finally, line 20 collects the query answer from all data silos.

To ease of understanding, Example 6 is used to further illustrate the main idea.

Example 6 Fig. 18 illustrates an instance as you mentioned, where k is 6 and the distance precision ϵ_0 is set to 1. Suppose the final upper bound u is 10, and the final lower bound l is 9.5. When using the upper bound u of 10 as the radius,

the number of points is 9, which exceeds k . When using the lower bound l of 9.5 as the radius, the number of points is 3, which is fewer than k . Since the difference between the bounds, $u - l = 0.5$, is less than the distance precision $\epsilon_0 = 1$, the binary search has been terminated and we can identify the case of ties. As shown in Fig. 18, this case introduces a challenge in handling boundary points to ensure that exactly k spatial objects are retrieved as the k nearest neighbors. To address this issue, we use the following strategy to manage these boundary spatial objects.

- (1) We perform a symmetric federated range query with the query object q (marked in red color in Fig. 18) and the circle radius $l = 9.5$. The query answer (*i.e.*, the three objects within the inner circle) will be cached individually within each data silo. Besides, these data silos collaboratively and securely compute the total count $k_l = 3$.
- (2) Now, it is still 3 objects short of the required $k = 6$ objects. Moreover, we do not know the exact number of ties in each data silo, but need to retrieve exactly 3 additional objects from tied ones. Thus, we need to sequentially request up to 3 points from each silo. Specifically, starting from data silo #1, we perform a symmetric federated range query with the query object q and the circle radius $u = 10$ in data silo #1, and limit the number of (uncached) query answer to 3. Although there are three objects within the query area, one of them has already been cached before, meaning data silo #1 can only provide two more answers as the k nearest neighbors. After this step, we need to request one more point from data silos #2 and #3. Similarly, we perform a symmetric federated range query with the query object q and the circle radius $u = 10$ in data silo #2, and limit the number of (uncached) query answer as 1. In data silo #2, three objects are within the query range, but one has been cached. Then, data silo #2 can randomly pick one from the two uncached points as the cached query answer. At this point, we have obtained exactly 6 spatial objects as the query answer, thus eliminating the need to request any additional (uncached) objects from data silo #3.
- (3) *Lastly*, a secure set union is performed to collect all the cached spatial objects from the data silos, resulting in a final output that contains exactly the 6 nearest neighbors.

Experimental Setup. Since we have not found such cases in our datasets (OSM and BJ), we generate a new synthetic dataset to conduct this experiment. In the new synthetic dataset, there are six data silos and each data silo has 10,000 spatial objects. We also fix k as 16 (*i.e.*, our default setting) and vary the proportion (denoted by α) of spatial objects strictly inside circle curve 25%, 50%, 75%, and 100%, respectively. We also ensure that at least 16 spatial objects happen to be on the circle curve (*i.e.*, they are “ties”). In other words, when $\alpha = 25\%$, $(1 - 25\%) \cdot k = 12$ spatial objects in the query answer are on the circle curve. Similarly, when $\alpha = 100\%$, no spatial object in the query answer is on the circle curve.

Experimental Result. Fig. 19 shows the experimental result. We can observe that our solution Hu-Fu is still more efficient than the other baselines. Additionally, we also evaluate the total cost of our solution to the above two technical issues (*i.e.*, identifying the case of ties and retrieving exactly k nearest neighbors). The total time cost of these two steps remains within 400ms, and the communication cost stays under 400KB, accounting for 4.4% of the total query time and 0.17% of the total communication cost.

Algorithm 4: Symmetric federated kNN query

Input: federation F , query object q , integer k
Output: the (exact) query answer ans

```

1  $[l, u] \leftarrow [0, U]$ , where  $U$  is the tight upper bound;
2  $q' \leftarrow$  secure location perturbation  $\text{Geol}(q)$ ;
3 while  $u - l \geq \epsilon_0$  do
4    $r \leftarrow (l + u)/2$ ,  $\mathcal{R} \leftarrow \text{circle}(q, r)$ ;
5    $v_i \leftarrow$  symmetric federated range counting in  $F$ 
     with circular range  $\mathcal{R}$ ;
6    $sign \leftarrow$  secure count comparison  $\text{CMP}(\{v_i\}, k)$ ;
7   if  $sign < 0$  then  $l \leftarrow r$ ;
8   else if  $sign > 0$  then  $u \leftarrow r$ ;
9   else break;
10 if  $sign = 0$  then
11    $ans \leftarrow$  symmetric federated range query in  $F$ 
     with circular range  $\text{circle}(q, r)$ ;
12 else // identify cases of ties & retrieve exact kNN
13    $\mathcal{R}_l \leftarrow \text{circle}(q, l)$ ,  $\mathcal{R}_r \leftarrow \text{circle}(q, r)$ ;
14    $k_l \leftarrow$  symmetric federated range counting in  $F$ 
     with circular range  $\mathcal{R}_l$  and cache all the spatial
     objects that are inside  $\mathcal{R}_l$  locally at each silo;
     // need to retrieve exactly remain objects
15    $remain \leftarrow k - k_l$ ;
16   foreach silo  $F_i \in F$  do // perform in sequential
17      $k_r^i \leftarrow$  symmetric federated range counting in
        $F_i$  with circular range  $\mathcal{R}_r$  over uncached
       spatial objects at silo  $F_i$ , limit the query
       answer to  $\leq remain$ , and cache the query
       answer locally at silo  $F_i$ ;
18      $remain \leftarrow remain - k_r^i$ ;
19     if  $remain = 0$  then break;
20    $ans \leftarrow$  secure set union  $\text{SUN}(S_1, \dots, S_n)$  over all
     cached spatial objects within each data silo;
21 return  $ans$ ;
```

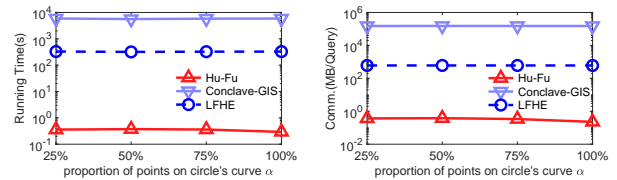


Fig. 19: Performance of symmetric federated kNN query in the synthetic dataset for handling ties.

D Additional Experiments on Heterogeneous Data Silos

Fig. 20 shows the results of asymmetric federated range query, distance join, and kNN join in the experiment of heterogeneous silos. We can first observe that our Hu-Fu can support all six spatial database systems, including PostGIS [10], MySQL [8], SpatiaLite [11], Simba [64], GeoMesa [27], and SpatialHadoop [21]. In addition, for all federated spatial queries in Fig. 20, we observe the running time of Hu-Fu becomes longer, compared with the results of homogeneous silos (*i.e.*, PostGIS in this test). This is because the time cost of Hu-Fu’s plaintext primitives (*i.e.*, plaintext range query) increases due to the slowest spatial database system. We can also observe that the experimental patterns of federated distance join and federated kNN join are similar to those of federated range query and federated kNN query. This is because a federated distance/kNN join is decomposed into a series of federated

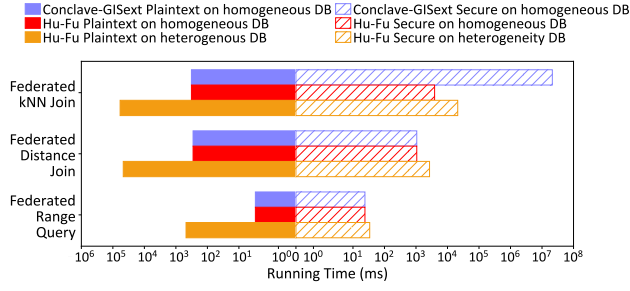


Fig. 20: Running time breakdown.

range/kNN query. Besides, we also recommend that all silos can use more efficient spatial database systems to further speed up the federated spatial queries.

This observation also holds for processing symmetric federated spatial queries, since the query processing procedure also involves plaintext primitives such as the plaintext range query and kNN query.

E Experiment on Hu-Fu with an Honest Broker

In the following, we conduct an experiment to demonstrate the performance of Hu-Fu with an honest broker (denoted by “Hu-Fu-HB”).

Experimental Setup. This experiment has tested all asymmetric federated spatial queries on the real-world dataset BJ. We compare Hu-Fu-HB with SMCQL-GIS/Conclave-GIS in query processing efficiency measured in running time and communication overhead. Besides, the other experiment settings are the same as those in Sec. 7.2.

Experimental Result. Fig. 21 and Fig. 22 present the experimental results of Hu-Fu-HB and SMCQL-GIS/Conclave-GIS in terms of running time and communication cost. First, We can observe that Hu-Fu-HB performs the same as SMCQL-GIS and Conclave-GIS on asymmetric federated range query and distance join. This is because all the compared algorithms (including ours) are simplified to Public due to the honest broker who collects the partial results (*i.e.*, sensitive data) of plaintext range query in each silo and is assumed to never reveal them to anyone else. Second, as for federated range counting, kNN query and kNN join, Hu-Fu-HB still outperforms SMCQL-GIS and Conclave-GIS with up to 4 orders of magnitudes faster in running time and 5 orders of magnitudes lower in communication cost. This is because our query writer is more effective to process queries on large-scale spatial data.

Summary. The experimental results demonstrate that compared to SMCQL-GIS and Conclave-GIS, Hu-Fu-HB exhibits similar efficiency in performing asymmetric federated range queries and distance joins. Additionally, Hu-Fu-HB achieves better efficiency in asymmetric federated range counting, kNN queries, and kNN joins. Although this experiment primarily focuses on asymmetric federated spatial queries, the observed patterns also apply to symmetric federated spatial queries. Specifically, the efficiency of symmetric federated range queries, distance joins, kNN queries, and kNN joins can be improved when an honest broker is present and a secure set union can be avoided. However, in contrast to these operations, symmetric federated range counting is not affected by the presence of an honest broker, as it does not involve a secure set union operator in its decomposition plan (as shown in Table 3).

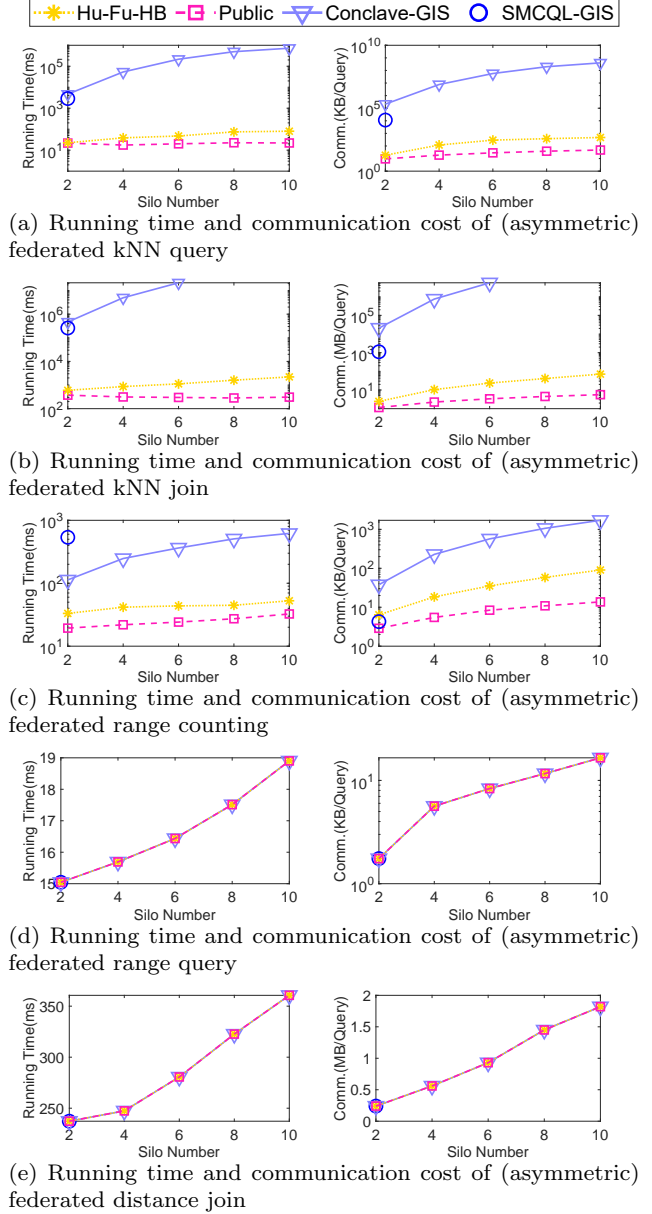


Fig. 21: Hu-Fu with an honest broker varying the silo number.

F Experiment on Geographically Partitioned Dataset

To explore the performance of Hu-Fu in the dataset where each silo corresponds to a single country, we have conducted an experiment on a new synthetic dataset (denoted by “OSM country”).

Experimental Setup. The OSM country dataset consists of spatial points from six countries in OpenStreetMap [9]. These countries include China, India, Japan, Malaysia, North Korea and South Korea, and each of them corresponds to a single data silo. As shown in Table 7, the volume of data in each silo follows the native data proportion of the corresponding country in OpenStreetMap. For example, the numbers of data

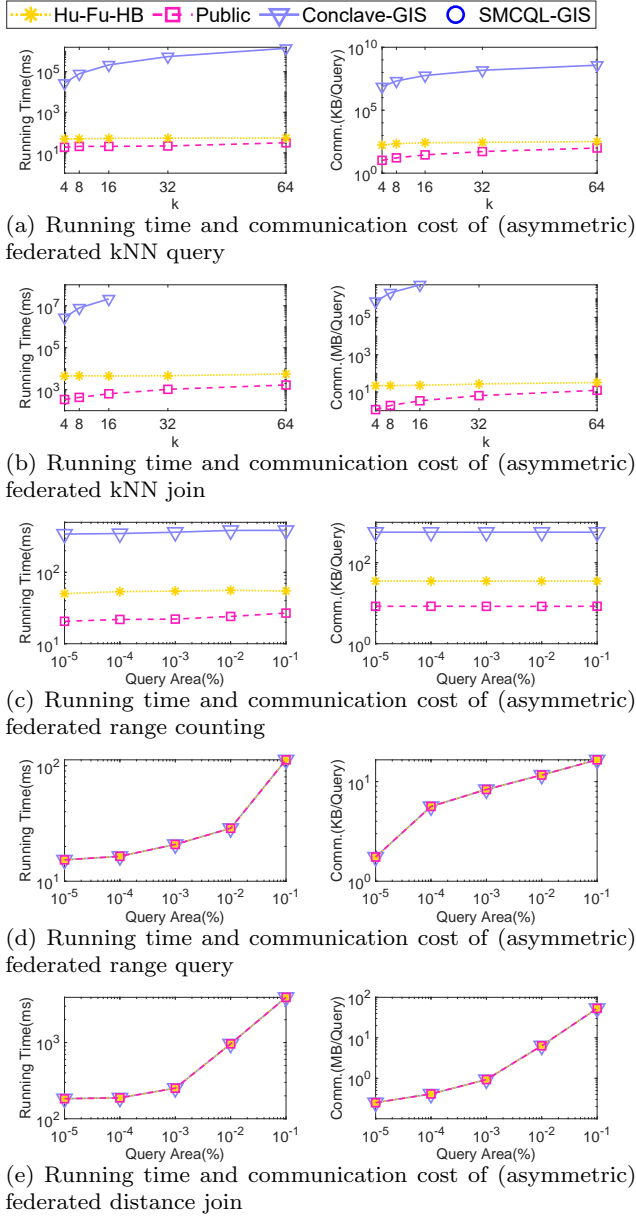


Fig. 22: Hu-Fu with an honest broker varying the query-specific parameter.

points in China, India and Japan are much larger than those of Malaysia, North Korea and South Korea. Here, we vary the total number of all data points from 10^4 to 10^9 . The query workloads are mainly asymmetric federated spatial queries. Besides, the other experiment parameters are the same as those in Sec. 7.4.1.

Experimental Result. The results of this experiment are shown in Fig. 23. First, we can observe that in the *OSM country* dataset, Hu-Fu also shows a significant improvement in both running time and communication cost over Conclave-GIS and Conclave-GISext on federated kNN query, kNN join and range counting. For example, Hu-Fu is at least 3 orders of magnitude faster than Conclave-GIS and Conclave-GISext in federated kNN query and kNN join, and costs 5 orders of magnitude lower network communication. Second, Hu-Fu

Table 7: Percentage of data of each country in *OSM country*.

Country	Percentage of data
China	21.7%
India	29.8%
Japan	39.4%
Malaysia	4.2%
North Korea	1.2%
South Korea	3.7%

achieves the same efficiency as Conclave-GISext on federated range query and distance join. Overall, the ranking of the compared solutions in terms of either running time or communication cost is consistent with the ranking in Sec. 7.4. Note that we exclude Conclave-GISext from the results of federated range counting as shown in Fig. 23c, because this query does not need the secure set union to protect data ownership in this query.

Summary. In this evaluation, Hu-Fu outperforms SMCQL-GISext and Conclave-GISext in terms of efficiency for asymmetric federated kNN queries, kNN join, and range counting. Meanwhile, Hu-Fu performs comparably to SMCQL-GISext and Conclave-GISext for asymmetric federated range queries and distance join. Thus, we can conclude that the experiment conducted on the *OSM country* dataset yields similar results to those observed in Sec. 7.4.1. This similarity may stem from the fact that none of the existing methods have leveraged the data distribution in their query processing strategies. Similarly, this new data partition is unlikely to alter the overall performance ranking of symmetric federated spatial queries.

G Experiment on the Improvement by the Index in Each Data Silo

To demonstrate that the federated spatial queries have already been accelerated by local indexes in Hu-Fu, we have conducted a new experiment on the efficiency improvement caused by these local indexes (*i.e.*, R-trees in our default setting).

Experimental Setup. In this experiment, we test the running time of the plaintext spatial queries in Hu-Fu (*i.e.*, plaintext range query, plaintext range counting and plaintext kNN query) on PostgreSQL with and without an R-tree. And the plaintext spatial queries are processed on the *OSM* dataset. Besides, we vary the data size from 10^4 to 10^8 and use the default setting of query area (0.001%) and k (16). The other experiment settings are the same as those in Sec. 7.1.

Experimental Result. As shown in Fig. 25, local indexes can improve the efficiency of plaintext spatial queries, especially when the data size is large. Specifically, the local index (R-tree) reduces the running time by up to $40\times$, $44\times$ and $2042\times$ when processing plaintext range query, range counting, and kNN query, respectively. Moreover, the plaintext spatial queries without local indexes cost even longer time than corresponding federated spatial queries in Hu-Fu. For instance, the running time of processing one federated kNN query by Hu-Fu takes 52 ms when the data size is 10^8 , which is even faster than that of the plaintext kNN query without the local index (2465 ms). Thus, the experimental result proves that we have used local indexes in Hu-Fu to speed up the processing of the plaintext operators.

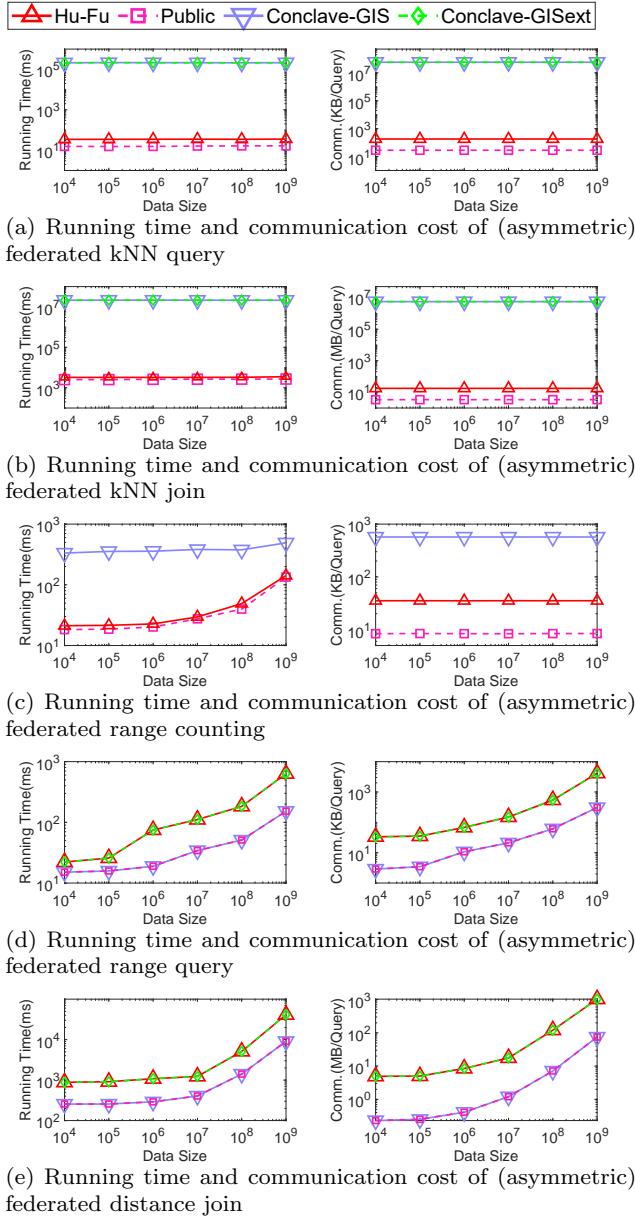


Fig. 23: Scalability test on synthetic dataset partitioned by countries.

H Impact of Additionally Protecting Query Privacy in Symmetric Queries on Time Cost

Symmetric queries generally exhibit slower performance compared to asymmetric queries. The gap stems from the essential difference in the problem definition (*i.e.*, whether query privacy needs to be protected or not). Consequently, accurately assessing the impact on time cost when incorporating query privacy protection is difficult.

However, we can roughly estimate the impact on time cost in two ways:

- By comparing the number of basic operators required for symmetric queries and asymmetric queries in Table 3 and Table 2, we can conclude that symmetric queries in-

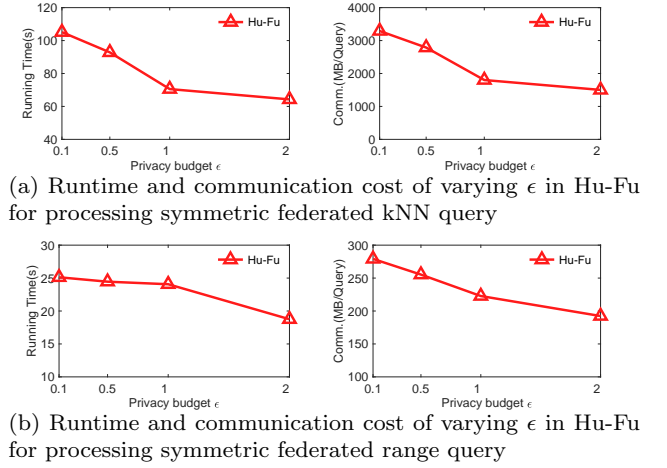


Fig. 24: Impact of privacy budget ϵ in Hu-Fu for processing symmetric federated kNN and range queries.

cur higher costs than asymmetric queries. For example, a symmetric federated range query involves many more secure operators (*e.g.*, secure distance comparisons) than an asymmetric federated range query. This explains why the former generally exhibits slower performance than the latter.

- We also report the impact on time cost under the default experimental setting. To obtain the result, we evaluate both Hu-Fu and Conclave-GIS for asymmetric federated range query and range counting on real-world dataset BJ in our new experimental environment. As shown in Table 8, a symmetric federated range query using Conclave-GIS is 4,054 \times slower than the corresponding asymmetric query. By contrast, our solution Hu-Fu significantly narrows this gap to 59 \times . We can observe a similar pattern for federated range counting. Although there is still a notable difference in the time cost between symmetric queries and asymmetric queries by using Hu-Fu, Table 8 demonstrates that Hu-Fu performs better in terms of query efficiency than the existing baseline, highlighting the challenge of achieving good efficiency for symmetric queries compared to asymmetric queries.

I Impact of Privacy Budget ϵ in Symmetric Queries

We use symmetric federated kNN and range queries to evaluate the impact of privacy budget ϵ on the query efficiency. The impact of ϵ for symmetric federated range counting and distance join is similar to that for symmetric federated range query, and the impact of ϵ for symmetric federated kNN join is similar to that of symmetric federated kNN query.

Specifically, we use the real-world dataset BJ to conduct the evaluation, and vary the privacy budget ϵ within a range from 0.1 to 2, where $\epsilon = 0.1$ indicates tighter privacy preservation level than $\epsilon = 2$. As for the other parameter settings, we use the default setting as introduced in Section 7.2.

Fig. 24 illustrates the running time and communication cost associated with symmetric federated kNN and range queries. As ϵ increases from 0.1 to 2, we observe a modest

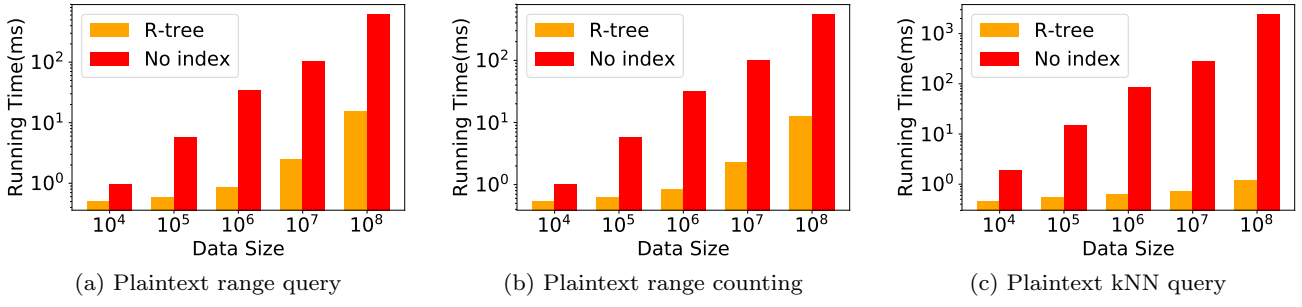


Fig. 25: Running time of plaintext spatial queries with/without R-tree.

Table 8: How many times slower is a symmetric query compared to an asymmetric query by specific solution

Solution	Federated Range Query	Federated Range Counting
Conclave-GIS	4054×	5302×
Hu-Fu	59×	77×

drop in both running time and communication cost for symmetric federated kNN query. A similar trend is also evident in the results of the symmetric federated range query. The overall pattern is reasonable, since a larger ϵ indicates a more relaxed privacy protection requirement, which results in a shorter distance between the query object and its perturbed location.