



武汉大学学报(信息科学版)

*Geomatics and Information Science of Wuhan University*

ISSN 1671-8860, CN 42-1676/TN

## 《武汉大学学报(信息科学版)》网络首发论文

题目： 面向大规模空间数据流的分布式连接查询方法  
作者： 李柞霖，向隆刚，余列冰，吴华意，关雪峰  
DOI： 10.13203/j.whugis20230040  
收稿日期： 2023-09-11  
网络首发日期： 2023-12-02  
引用格式： 李柞霖，向隆刚，余列冰，吴华意，关雪峰. 面向大规模空间数据流的分布式连接查询方法[J/OL]. 武汉大学学报(信息科学版).  
<https://doi.org/10.13203/j.whugis20230040>



**网络首发：**在编辑部工作流程中，稿件从录用到出版要经历录用定稿、排版定稿、整期汇编定稿等阶段。录用定稿指内容已经确定，且通过同行评议、主编终审同意刊用的稿件。排版定稿指录用定稿按照期刊特定版式（包括网络呈现版式）排版后的稿件，可暂不确定出版年、卷、期和页码。整期汇编定稿指出版年、卷、期、页码均已确定的印刷或数字出版的整期汇编稿件。录用定稿网络首发稿件内容必须符合《出版管理条例》和《期刊出版管理规定》的有关规定；学术研究成果具有创新性、科学性和先进性，符合编辑部对刊文的录用要求，不存在学术不端行为及其他侵权行为；稿件内容应基本符合国家有关书刊编辑、出版的技术标准，正确使用和统一规范语言文字、符号、数字、外文字母、法定计量单位及地图标注等。为确保录用定稿网络首发的严肃性，录用定稿一经发布，不得修改论文题目、作者、机构名称和学术内容，只可基于编辑规范进行少量文字的修改。

**出版确认：**纸质期刊编辑部通过与《中国学术期刊（光盘版）》电子杂志社有限公司签约，在《中国学术期刊（网络版）》出版传播平台上创办与纸质期刊内容一致的网络版，以单篇或整期出版形式，在印刷出版之前刊发论文的录用定稿、排版定稿、整期汇编定稿。因为《中国学术期刊（网络版）》是国家新闻出版广电总局批准的网络连续型出版物（ISSN 2096-4188，CN 11-6037/Z），所以签约期刊的网络版上网络首发论文视为正式出版。

Doi: 10.13203/j.whugis20230040

### 引用格式:

李柞霖, 向隆刚, 余列冰, 等. 面向大规模空间数据流的分布式连接查询方法[J]. 武汉大学学报(信息科学版), 2023, Doi:10.13203/j.whugis20230040. (LI Zuolin, XIANG Longgang, YU Liebing, et al. Distributed Join Query Method for Large-Scale Spatial Data Streams[J] *Geomatics and Information Science of Wuhan University*, 2023, Doi:10.13203/j.whugis20230040.)

## 面向大规模空间数据流的分布式连接查询方法

李柞霖<sup>1</sup> 向隆刚<sup>1</sup> 余列冰<sup>1</sup> 吴华意<sup>1</sup> 关雪峰<sup>1</sup>

1 武汉大学测绘遥感信息工程国家重点实验室, 湖北 武汉, 430079

**摘要:**空间连接查询是处理和分析空间数据的基础操作之一。随着空间数据的爆发式增长, 针对海量空间数据的连接查询技术备受瞩目。面向大规模历史数据的空间连接查询已被广泛研究, 而受限于数据流的高速接入率与连接的实时性需求, 目前面向数据流的分布式空间连接查询仍充满挑战。为此, 本文面向大规模空间数据流设计了一种分布式连接查询处理框架。首先形式化定义了针对空间数据流连接查询处理问题, 按照参与连接的数据形态细分为“流-表”和“流-流”两类连接。进一步分析分布式连接场景的共性问题, 设计了“全局网格分区-局部空间索引”的两层处理框架, 以支持空间流的分布式连接。在此基础上, 针对不同连接场景分别设计了适应的连接策略: 对于“流-表”连接, 提出了基于两级 R-tree 拓扑关系判断优化算法; 对于“流-流”连接, 设计了一种顾及分区边界的数据冗余路由算法, 以保证分区边界数据的正确连接; 此外, 针对间隔时间语义的缓存需求, 提出了兼顾状态管理与检索效率的 BinR-tree 结构。大量实验结果表明, 本文提出的空间数据连接方法具有良好的线性加速比, 且相对于基线方法, 连接查询效率得到了显著提升。

**关键词:**空间数据流; 空间连接; 流连接; 分布式; 空间索引

## Distributed Join Query Method for Large-Scale Spatial Data Streams

LI Zuolin<sup>1</sup> XIANG Longgang<sup>1</sup> YU Liebing<sup>1</sup> WU Huayi<sup>1</sup> GUAN Xuefeng<sup>1</sup>

1 State Key Laboratory of Information Engineering in Surveying Mapping and Remote Sensing, Wuhan University, Wuhan 430079, China

**Abstract: Objectives:** Spatial join query is one of the basic operations for processing and analyzing spatial data. With the explosive growth of spatial data, join query technology for massive spatial data has attracted much attention. Although the existing research work has explored the join query of spatial data streams, it is still in its infancy. The processing of join query processing of spatial data streams is still insufficient in terms of systematization and universality. It is urgent to explore how to define general Spatial data streams connect query problems and provide systematic solutions. **Method:** After fully excavating the common real-time spatial data processing problems behind various practical application scenarios, this paper refers to the processing theory and methods of spatial data in the batch processing field, and considers the characteristics of long-term continuous operation of stream processing, and formalizes the problem of spatial data stream connection Two types of connection operators are defined, including "stream-table" connection and "stream-stream" connection; on this basis, a two-layer data index structure of global grid partition and local spatial index is proposed to support spatial Stream supports distributed connection; for the "stream-table" connection, this paper proposes two physical realization methods of spatial dimension tables, and designs a two-level R-tree-

收稿日期:2023-09-11

项目资助: 湖北珞珈实验室专项基金 (220100010)。

第一作者:李柞霖, 硕士。lynnlee1229@whu.edu.cn

通讯作者:向隆刚, 博士, 教授, 主要从事空间数据库、轨迹数据时空挖掘研究。geoxlg@whu.edu.cn

based topology relationship judgment optimization for the memory-based storage method Algorithm; For "stream-to-stream" joins, this paper proposes a partition boundary data redundancy algorithm to correctly implement cross-partition joins of partition boundary data. In addition, for the caching requirements of interval time semantics, a BinR-tree result that takes into account both management and retrieval efficiency is proposed. **Results:** (1) The "stream-table" spatial connection implemented in this paper can achieve a throughput of more than 60,000 under a single degree of parallelism, and the average overall delay is about 90 milliseconds. Compared with the native method, the average throughput is 9 times Boost, latency is reduced by an average of 20 milliseconds. (2) Same as the "stream-table" connection, the performance of the "stream-stream" connection method in this paper improves steadily with the increase of parallelism, and has good horizontal scalability. (3) In the experimental comparison between window connection and interval connection, with As the window increases, the connection latency increases, much higher than the stream interval connection. This part of the experiment also shows that the stream interval connection without bounded time semantics is more suitable for stream processing logic. (4) Mesh division has little effect on throughput, and its main function is to balance the load and ensure the spatial proximity of data. (5) When the bin size is smaller than the window size, the query will span multiple BinR-trees, which will affect the query efficiency. Properly increasing the bin size will improve the query efficiency; when the bin size is greater than or equal to the time interval, the query efficiency will not be greatly improved. **Conclusions:** A large number of experimental results show that the spatial data connection method proposed in this paper has a good linear speedup ratio, and compared with the baseline method, the connection query efficiency has been significantly improved.

**Key words:** spatial data streams; spatial join; stream join; distributed; spatial index

随着传感器技术的发展、物联网的建设以及云计算的普及,海量的空间数据以数据流的形式源源不断地传送到数据中心,交通出行、移动社交等越来越多地理信息相关的行业需要对这些数据流进行持续的实时处理,从而为复杂的时空问题提供低延迟的解决方案。近年来,随着空间数据的爆发式增长,对海量数据的查询处理技术备受瞩目<sup>[1-2]</sup>。空间连接查询是处理和分析空间数据的基础操作之一,在如交通流量分析、区域监控、人车匹配等诸多场景中具有重要应用价值<sup>[2-6]</sup>。

**定义 1.** 空间数据流。空间数据流表示无界的空间数据集合,即空间数据可以持续生成,无限增长,用以描述空间对象的实时状态,可以形式化定义为由空间对象  $o$  及其数据生产时间戳  $t$  组成的二元组  $\langle o, t \rangle$  的集合。

$$S = \{ \langle o_i, t_i \rangle | i = 1, 2, \dots \} \quad (1)$$

针对空间数据流的连接查询处理受到了越来越多的关注。与针对大规模历史数据表的空间连接查询不同,面向空间数据流的连接查询更具有挑战。传统的数据库大都采用索引嵌套循环连接 (indexed nested loop join, INLJ) 算法<sup>[7]</sup>实现两个表的空间连接。然而,空间数据流高速接入且无限增长的特点,使其难以直接使用 INLJ 算法;在流场景下,连接中某个流中的记录到达时,另一个流中满足连接条件的记录可能尚未到达,因而必须记录数据流的状态,从而在适当的时机触发计算;此外,还必须维护和清理数据流的过期状态,防止其无限增长。

面向大规模历史空间数据的连接查询处理技术伴随着分布式批处理引擎 (distributed batch processing

engine, DBPE) 的不断发展而日趋成熟。依据所依赖的底层计算引擎,基于 DBPE 的空间连接查询实现主要分为两类。一类是基于 MapReduce 计算引擎的扩展,如 Hadoop-GIS<sup>[8]</sup>、SpatialHadoop<sup>[9]</sup>等。其中, Hadoop-GIS 是 Apache Hadoop 之上第一个关于空间数据处理的扩展,它将 Apache Hadoop 视为黑盒并完全依赖其底层架构,在此基础上提供了基于均匀网格的数据分区策略,在各个分区内支持构建 R\*-tree 索引加速计算; SpatialHadoop 对 Apache Hadoop 的源码进行了扩展,使其从存储和计算两个方面原生支持空间数据,从而获得了比 Hadoop-GIS 更好的性能。另一类是基于 Apache Spark 的扩展,如 SpatialSpark<sup>[10]</sup>、GeoSpark<sup>[11-12]</sup>、Simba<sup>[13]</sup>、LocationSpark<sup>[14]</sup>等。SpatialSpark 和 GeoSpark 在 Spark 之上提供了空间数据分区策略,并支持多种类型的空间查询。而 Simba 和 LocationSpark 则更关注对 SQL 的扩展。现有研究大多支持点、线、矩形以及多边形等多种空间数据类型,并在此基础上支持多种类型的空间连接,如距离连接 (distance join)、拓扑连接 (topology join)、k-最近邻连接 (k nearest neighbors join, kNN Join) 等<sup>[15-16]</sup>。然而,受限于 DBPE 的执行模式,它们的高吞吐以高延迟为代价。对于持续接入的空间数据,这些框架难以进行低延迟的连接查询处理。

受到历史空间数据连接查询处理的启发,以及分布式流处理引擎 (distributed streaming processing engine, DSPE) 的日臻成熟,空间数据流的连接查询处理技术也得到了一定的发展。Zhang<sup>[17]</sup>率先基于 Apache Storm 实现了轨迹数据流的距离连接,这一工作在空间数据流的

低延迟连接查询处理方面具有重要的借鉴意义，但其仅支持基于轨迹点距离的连接；MobyDick<sup>[18]</sup>构建在Apache Flink之上提供了空间流连接支持，但其研究侧重于对空间数据流处理提供类SQL支持，缺少针对性的分区和索引策略；GeoFlink<sup>[19-20]</sup>同样基于Apache Flink，支持空间数据流的距离连接，不过其完全采用有界空间数据的查询语义，在流处理场景下使用范围有限。

尽管现有研究工作已经对空间数据流的连接查询进行了一定的探索，但是仍处于起步阶段，空间数据流连接查询处理在系统性和普适性方面仍存在不足，亟需探索如何定义通用的空间数据流连接查询问题，并提供系统化的解决方案。为此，本文在充分挖掘各类实际应用场景背后共性的即时空间数据处理问题后，借鉴批处理领域空间数据的处理理论与方法，考虑流处理长期持续运行的特点，针对空间数据流连接问题形式化定义了2类连接算子，包括“流-表”连接和“流-流”连接；在此基础上，提出了一种全局网格分区及局部空间索引的两层数据索引结构，以支持空间数据流分布式连接；对于“流-表”连接，本文提出了两种空间维度表的物理实现方式，并针对其中基于内存的存储方式设计了一种基于两级R-tree的拓扑关系判断优化算法；对于“流-流”连接，本文提出了一种分区边界数据冗余算法，以正确实现分区边界数据的跨分区连接。此外，针对间隔时间语义的缓存需求，提出了兼顾管理与检索效率的BinR-tree结构。

## 1 问题定义

本文研究面向空间数据流的连接查询方法，首先给出本文所指空间连接查询的形式化定义。

**定义 2.** 空间连接查询。给定空间数据集  $O_1, O_2$ ，对于  $\forall o_1 \in O_1$ ，找到所有  $o_2 \in O_2$ ，满足  $r(o_1, o_2)$  并输出连接结果  $j(o_1, o_2)$ 。其中，函数  $r$  用于判断两个对象之间的

空间关系是否符合连接条件，除支持九交模型所定义的所有拓扑关系<sup>[21]</sup>外，本文增加了距离关系（包括距离包含于、距离包含），相应地，定义  $r_\delta$  表示考虑距离阈值  $\delta$  的空间关系判断函数。连接函数  $j$  用于对连接对象进行定制化的合并操作。空间连接的形式化定义如式(2)。

$$Join(O_1, O_2, \delta) = \{j(o_1, o_2) \mid o_1 \in O_1, o_2 \in O_2, r_\delta(o_1, o_2) = true\} \quad (2)$$

$$r_\delta(o_1, o_2) = \begin{cases} r(o_1, o_2) \wedge dis(o_1, o_2) < \delta & \delta > 0, \\ r(o_1, o_2) & \delta \leq 0 \end{cases} \quad (3)$$

式中，dis 函数用于计算空间对象间距离。

根据参与连接对象  $O_1, O_2$  的数据形态分类，本文基于实际应用场景抽象了两类空间数据流连接问题，包括用于空间数据流与空间维度表的“流-表”连接，以及用

于“流-流”的空间窗口连接和空间间隔连接。其中，“流-表”连接用于将单空间数据流关联静态空间表信息，例如将用户移动轨迹与行政区划连接以追踪其所经过的行政单元；“流-流”连接适用于连接对象均为空间数据流的场景，例如在交通领域中常需要把行人轨迹数据流与网约车轨迹数据流进行空间连接，从而为行人确定附近的网约车。下面给出相关的形式化定义。

### 1.1 “流-表”空间连接

维度数据是由独立的、不重叠的数据元素组成的数据集，主要有过滤、分组和标签三个方面的功能<sup>[22]</sup>。在空间数据流处理中，同样有为流几何关联维度信息的需求。为此，本文将维度数据的概念引入空间数据处理领域，将带有空间属性的维度数据称为空间维度数据，其形式化定义如下。

**定义 3.** 空间维度数据。空间维度数据是由独立的、不重叠的且带有空间属性的数据元素组成的数据集。空间维度数据中的单个数据元素称为维度几何。

基于空间维度数据的定义，本文把空间数据流与空间维度表基于拓扑关系的实时连接称为“流-表”空间连接，其形式化定义如下。

**定义 4.** “流-表”空间连接。给定空间数据流  $S$ ，空间维度表  $D$ ，对  $\forall s \in S$ ，找到所有  $d \in D$ ，满足  $r_\delta(s, d)$  并输出连接结果  $j(s, d)$ 。

### 1.2 “流-流”空间连接

#### 1.2.1 “流-流”空间窗口连接

本文所定义的“流-流”空间窗口连接面向两个高速变化的空间数据流，将落在同一时间窗口内的两个窗口快照按照指定拓扑关系进行连接，其相关形式化定义如下。

**定义 5.** 窗口快照。如图 1，给定空间数据流  $S$ ，起始时间  $T$ ，窗口大小  $W$ ，窗口滑动  $\varepsilon$ ， $S$  中落在任意时间窗口  $[T + i \cdot \varepsilon, T + i \cdot \varepsilon + W]$  范围内的所有空间数据形成空间数据流  $S$  的一个窗口快照。

**定义 6.** “流-流”空间窗口连接。给定两个空间数据流  $S_1$  和  $S_2$ ，对于分别属于  $S_1$  和  $S_2$  且落在同一窗口内的窗口快照  $S_{1w}$  和  $S_{2w}$ ，“流-流”空间窗口连接是对  $\forall s_1 \in S_{1w}$ ，找到所有  $s_2 \in S_{2w}$ ，满足  $r_\delta(s_1, s_2)$ ，输出  $j(s_1, s_2)$ 。



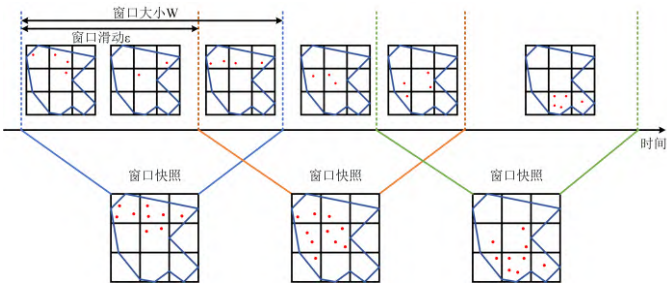


图 1 窗口快照示意图  
Fig. 1 Window Snapshot

### 1.2.2 “流-流”空间间隔连接

基于窗口的空间流连接，只能对当前窗口内数据进行连接，无法处理相对向前、向后时间数据。为此，本文定义“流-流”空间间隔连接以拓展空间流连接的时间语义，其形式化定义如下。

**定义 6.** “流-流”空间间隔连接。给定两个空间数据流  $S_1$  和  $S_2$ ，上界时间间隔  $T_{ub}$  和下界时间间隔  $T_{lb}$ 。对  $\forall s_1 \in S_1, \forall s_2 \in S_2$  设其时间戳分别为  $t_1, t_2$ 。“流-流”空间间隔连接是对  $\forall s_1 \in S_1$ ，找到所有  $s_2 \in S_2$ ，且  $t_1 + T_{lb} \leq t_2 \leq t_1 + T_{ub}, r_\delta(s_1, s_2)$ ，输出  $j(s_1, s_2)$ 。

为方便描述，下文称“流-流”连接中， $S_1$  为空间驱动流， $S_2$  为空间被驱动流。

## 2 空间数据流连接查询算法

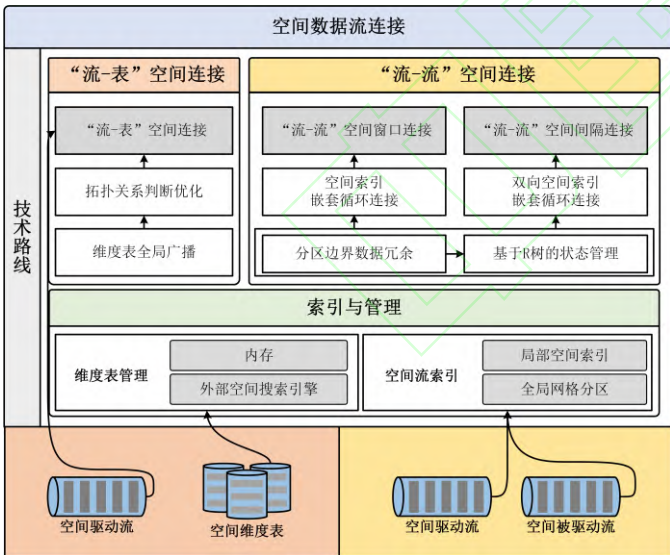


图 2 空间数据流连接技术路线

Fig. 2 Technical Routes of the Spatial Data Streams Join

本文研究开展的技术路线如图 2。首先，对于输入空间数据进行适应分布式连接的索引与管理，对空间流数据建立“全局网格分区-局部空间索引”两级索引，对空间维表根据其物理存储形式进行对应的管理。在此基础上，对各连接算子进行针对性优化。“流-表”空间连

接基于广播表的思想实现，并进一步对拓扑判断进行优化；对于“流-流”空间连接，本文顾及数据分区、状态管理等问题，将嵌套循环连接算法推广至空间流处理。

本文提出的查询方法针对空间数据流和分布式环境的特点进行了分析，设计了适用于空间流连接算子的共用索引与管理框架。在此基础上，针对不同算子进行了适应性的调整与进一步优化。

### 2.1 空间流索引

面向大规模空间数据的分布式连接查询实现一般包含两层索引。第一层用于全局的数据分区以便在分布式环境下对数据进行划分，这一层索引既要实现数据的均匀划分，又要保持每个节点内数据的空间邻近性，以尽可能减少处理时各个计算节点间的通信成本；第二层作为局部索引在每个节点上加速计算。空间数据索引可大致分为两类。1) 网格状索引，如均匀网格 (uniform grid)、Z-curve、Hilbert 等；2) 树状索引，如 R-tree、Quadtree 等<sup>[3][23]</sup>。

相较于网格状索引，树状索引的查询效率更高，基于 DBPE 的一般采用两层树状索引的结构。然而，树状索引的维护成本较高，持续到达的空间数据流可能带来频繁的索引重构。网格状索引在给定数据边界的情况下结构固定且不需要维护。因此，本文选择使用基于规则格网的全局分区索引，在各分布式节点使用局部树状索引加速空间查询。

### 2.2 “流-表”空间连接

“流-表”空间连接的计算图如图 3 所示，空间维度表以广播的形式发往每个连接算子，故无需对空间数据流构建索引。连接算子读取来自多个输入分区的数据流，根据输入条件立即对每个到达的几何对象进行判断，符合连接条件的几何对象与对应维度表连接。

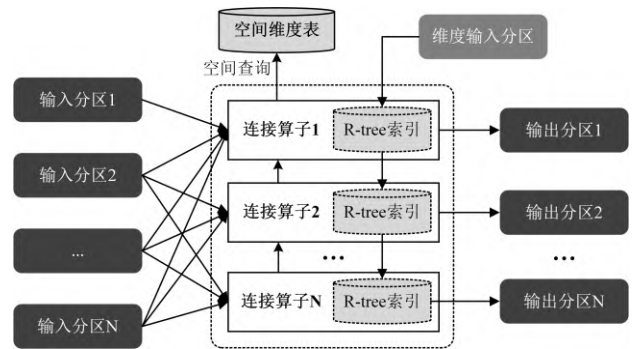


图 3 “流-表”空间连接计算图

Fig. 3 Stream-Table Spatial Join Calculation Graph

连接算子在初始化阶段可选择是否为本文提供两种空间维度表的物理实现方式：其一是内存存储，将维度表作为连接算子的本地状态。另一种方式是把空间检索和拓扑关系判断全权交给支持高并发、低延迟空间检索的外部引擎，如 GeoMesa<sup>[24]</sup>。虽然在实现上相对简单，

但是由于需要大量访问外部存储，在效率上不如前者。

在流表连接流程中，最关键的问题是连接算子如何高效判断流几何与查询几何之间的拓扑关系。为此，本文提供了两个层面的优化，第一个层面是在查询几何数量较多的情况下对其建立 R-tree 索引，以加速空间搜索效率；第二个层面是对常见的拓扑关系判断进行了特殊的优化。本文以最为常见的点与多边形之间的包含关系的判断为例进行优化介绍。

射线法常用于判断点是否包含于多边形内部<sup>[25]</sup>。现实场景中常存在复杂多边形，如未经简化的行政区划多边形，其边数通常在数千到数万不等。针对此类复杂多边形，本文提出了一种基于 R-tree 优化的射线法。该算法将多边形的各条边视为独立的几何对象，并为其构建 R-tree 索引。

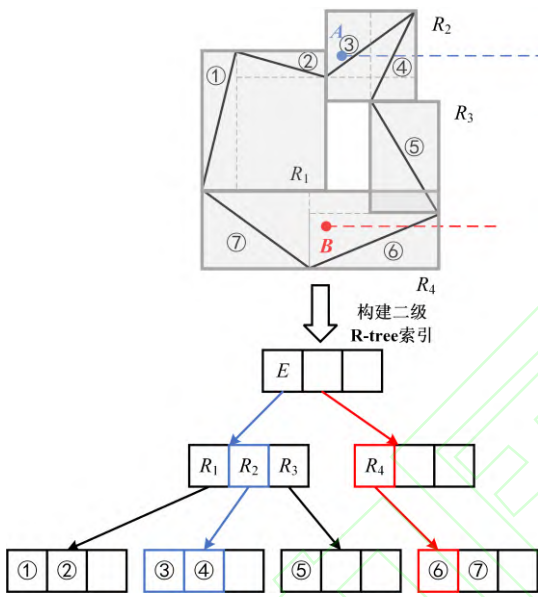


图 4 基于 R-tree 的射线法

Fig. 4 R-tree based ray method

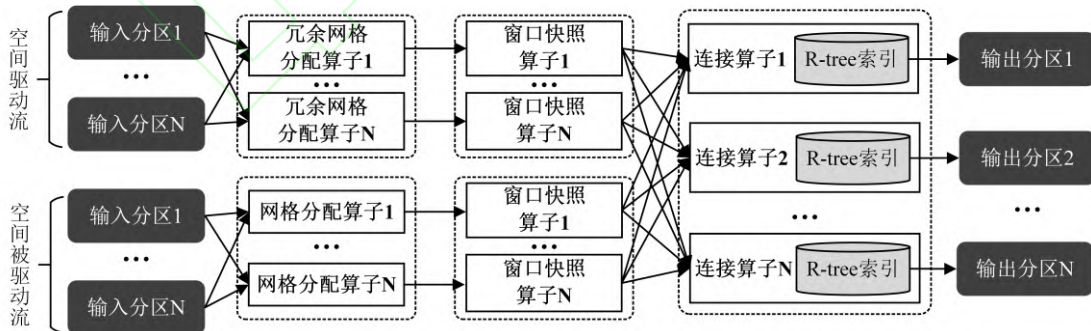


图 5 “流-流”空间窗口连接计算图

Fig. 5 Stream-Stream Spatial Window Join Calculation Graph

SINLJ 算法的流程如算法 2 所示，首先为所有来自被驱动流的几何对象创建 R-tree 索引，然后遍历所有来自驱动流的几何对象，在遍历过程中通过 R-tree 索引检索出被驱动流中与之匹配的候选几何，最后判断与候选几何之间的拓扑关系是否符合需求，将满足要求的对象

如图 4，A、B 分别为点在多边形外部和内部的情况，传统射线法需要将以 A、B 为起点的射线与多边形所有边进行相交判断，而借助索引可快速定位可能与射线相交的边。以点 A 为例，首先在 R-tree 中查询到其包围盒为  $R_2$ ，则只需将以 A 为起点的射线与  $R_2$  中所有可能相交边（③、④）进行相交判断即可确定该点与多边形的位置关系，该方法可将时间复杂度从  $O(n)$  降低到  $O(\log n)$ 。基于 R-tree 优化的射线法的算法流程如算法 1 所示。

算法 1. 基于 R-tree 优化的射线法

输入：空间点  $P$ ，多边形  $R$ ；  
输出：点是否在多边形内判读结果  $result$ 。

- ① R-tree  $rt \leftarrow \emptyset$
- ② for each edge  $e$  in  $R$  do
- ③  $rt.insert(e)$
- ④  $line \leftarrow Ray(P)$
- ⑤  $result \leftarrow false$
- ⑥ for each edge  $e$  in  $rt.query(line)$  do
- ⑦ if  $Cross(e, line)$  do
- ⑧  $result \leftarrow !result$

### 2.3 “流-表”空间连接

#### 2.3.1 “流-流”空间窗口连接

“流-流”空间窗口连接的计算图如图 5 所示。网格分配算子从输入分区中读取几何对象并为其生成网格索引，之后根据网格索引对流几何进行重新分发，根据用户指定的窗口大小划分为窗口快照后进行连接处理。本文借鉴了 INLJ 算法，将其调整后用于空间数据的连接，本文称之为空间索引嵌套循环连接（spatial index nested loop join, SINLJ）算法。

连接后输出。

算法 2. 空间索引嵌套循环连接算法

输入：窗口快照  $S_{1w}$ ,  $S_{2w}$ ，空间关系判断函数  $r\delta$ ；  
输出：符合连接条件的子对  $j(s_1, s_2)$ 。

- ① R-tree  $rt \leftarrow \emptyset$



- ② for each  $s_2$  in  $S_{2w}$  do
- ③      $rt.insert(s_2)$
- ④ for each  $s_1$  in  $S_{1w}$  do
- ⑤     for each  $s_2$  in  $rt.query(s_1)$  do
- ⑥         if  $r_\delta(s_1, s_2)$  do
- ⑦             output  $j(s_1, s_2)$

在空间连接中，落在不同网格内的几何对象很可能出现跨分区连接的现象。为了解决这一问题，本文设计了一种分区边界数据冗余算法。

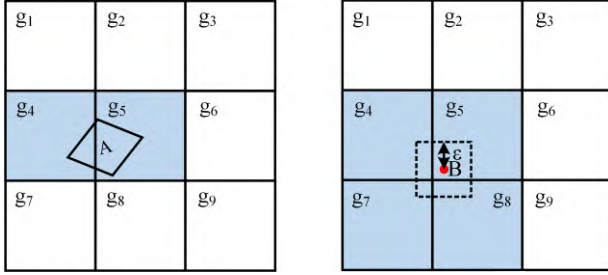


图 6 驱动流几何网格分配案例

Fig. 6 Driving Stream Geometry Mesh Assignment Example

具体而言，对于驱动流中的几何对象，若其为非点几何对象则为其分配所覆盖的所有分区编号，如图 6 中的多边形  $A$  所对应的网格编号为  $g_4$  和  $g_5$ ；若其为点对象且为距离连接，那么为其分区  $\varepsilon$  领域所覆盖的所有分区编号，如图 6 中的点  $B$  所对应的网格编号为  $g_4$ 、 $g_5$ 、 $g_7$  和  $g_8$ 。分区边界冗余网格分配算法的具体流程如算法 3 所示。

算法 3. 分区边界冗余网格分配算法

输入：几何对象  $P$ ，距离阈值  $\varepsilon$ （仅用于点的距离连接），网格边长  $l$ ；

输出：网格编号  $(i, j)$ 。

- ①  $box \leftarrow GetBoundingBox(P, \varepsilon)$
- ②  $\min X \leftarrow box.min Lng / l$
- ③  $\max X \leftarrow box.max Lng / l$
- ④  $\min Y \leftarrow box.min Lat / l$
- ⑤  $\max Y \leftarrow box.max Lat / l$
- ⑥ for  $i$  from  $\min X$  to  $\max X$  do
- ⑦     for  $j$  from  $\min Y$  to  $\max Y$  do
- ⑧         output  $(i, j)$

需要注意的是，在“流-流”空间窗口连接中，仅对驱动流使用分区边界冗余的网格分配算法，而被驱动流采用常规的网格分配算法，否则将导致出现重复的连接结果。另外，由于冗余网格分配算法的使用，“流-流”空

间窗口连接仅支持驱动流包含非点几何对象。为了使连接更加高效，在创建作业时推荐将数据量较小的空间数据流作为驱动流。

图 7 是一个基于 SINLJ 算法的“流-流”空间窗口连接示例，在该示例中两个空间数据流均由空间点组成，窗口快照内点的空间分布如图右上角矩形虚线框内所示。依据 SINLJ 算法，窗口快照内属于被驱动流的点将被插入到 R-tree 索引中。对于窗口快照内属于驱动流的每个点，通过 R-tree 索引检索出候选的点，图中仅以空间点  $P$  为例，可在 R-tree 索引中通过路径  $R_5 \rightarrow R_2$ ，快速确定与  $P$  关联的候选几何对象④、⑤，进一步距离判断⑤符合连接条件，连接输出。

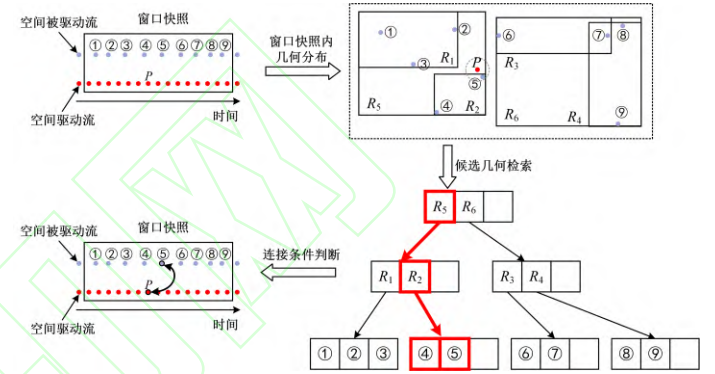


图 7 “流-流”空间窗口连接示例

Fig. 7 Stream-Stream Spatial Window Join Example

### 2.3.2 “流-流”空间间隔连接

“流-流”空间间隔连接的计算图如图 8 所示，相比于“流-流”窗口连接，舍弃了有界的窗口时间语义，在连接处理时更加灵活，但也带来了新的问题。每个空间数据流中的几何对象到达时，均无法确定另一个数据集中落在指定时间范围内的数据是否已经全部到达，难以直接采用 SINLJ 算法进行连接。为此，本文在 SINLJ 的基础之上提出了双向空间索引嵌套循环连接（bidirectional spatial index nested loop join, BSINLJ）算法。

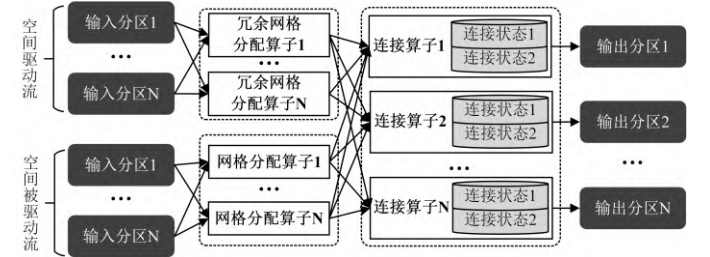


图 8 “流-流”空间间隔连接计算图

Fig. 8 Stream-Stream Spatial Interval Join Calculation Graph

BSINLJ 算法的流程如算法 4 所示，首先创建两个状态  $st_1$  与  $st_2$  分别用于存储来自两个空间数据流的几何对象，对于每个到来的属于  $S_1$  的几何对象  $s_1g_1$ ，检索  $st_2$  中  $[Ts_1 + T_{lb}, Ts_1 + T_{ub}]$  时间间隔内的候选数据，进一步筛选

出符合空间连接条件的几何对象，连接输出，最后将  $s_1$  插入  $st_1$  中以记录状态。 $S_2$  中的几何对象作对称的操作。

**算法 4.** 双向空间索引嵌套循环连接算法

输入：空间数据流  $S_1$ ,  $S_2$ , 上界时间  $T_{ub}$ , 下界时间  $T_{lb}$ , 空间关系判断函数  $r_\delta$ ;

输出：符合连接条件的子对  $j(s_1, s_2)$ .

- ① State  $st_1 \leftarrow \emptyset$
- ② State  $st_2 \leftarrow \emptyset$
- ③ for each  $s_1$  in  $S_1$  do
- ④     for each  $s_2$  in  $st_2.query(s_1)$  do
- ⑤         if  $Ts_1 + T_{lb} \leq T(s_2) \leq Ts_1 + T_{ub}$  and  $r_\delta(s_1, s_2)$  do
- ⑥             output  $j(s_1, s_2)$
- ⑦          $st_1.insert(s_1)$
- ⑧ for each  $s_2$  in  $S_2$  do
- ⑨     for each  $s_1$  in  $st_1.query(s_2)$  do
- ⑩         if  $Ts_2 - T_{ub} \leq T(s_1) \leq Ts_2 - T_{lb}$  and  $r_\delta(s_1, s_2)$  do
- ⑪             output  $j(s_1, s_2)$
- ⑫      $st_2.insert(s_2)$

BSINLJ 算法中状态的管理对于其正确执行至关重要。

<时间戳, 数据列表>是常用的存储结构, 只需要判断时间戳是否过期。然而, 该结构查询时需要对所有数据遍历计算, 对于计算开销大的空间数据有较大优化空间。沿用 R-tree 是自然的优化思路, 但是单 R-tree 结构数据

删除较耗时, <时间戳, R-tree>结构又面临 R-tree 频繁建立的问题。为此, 本文基于数据分箱 (binning) 思想提出了 BinR-tree 状态管理方法, 即对于空间数据流按照其时间戳分箱后, 再建立 R-tree 索引, 将数据按照 <BinID, R-tree>形式存储。如图 9。该方法既保留了主流存储结构状态管理高效的优点, 又提升了空间数据查询效率。

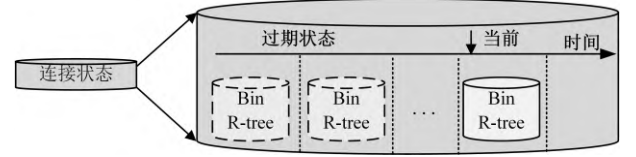


图 9 BinR-tree 结构

Fig. 9 Structure of BinR-tree

对应地, 算法 4 中的查询和插入操作都需要先进行计算 BinID, 计算公式如下:

$$BinID = \frac{timestamp}{Binsize} \quad (4)$$

连接状态存储于内存, 随着时间的推移, 过期状态的堆积不仅会影响查询性能, 一旦超过限制还会造成内存溢出错误。本文的状态清理方案为在建立 BinR-tree 的同时设定其对应的清除触发器, 按照式(5)计算触发时间:

$$\max\{t_{binmax}, t_{binmax} + T_{rub}\} \quad (5)$$

其中,  $t_{binmax}$  是箱的时间上界,  $T_{rub}$  指相对上界时间。对于  $S_1$  中元素, 相对时间上界为  $T_{ub}$ 。对于  $S_2$  中元素, 其相对时间上界为  $-T_{lb}$ 。

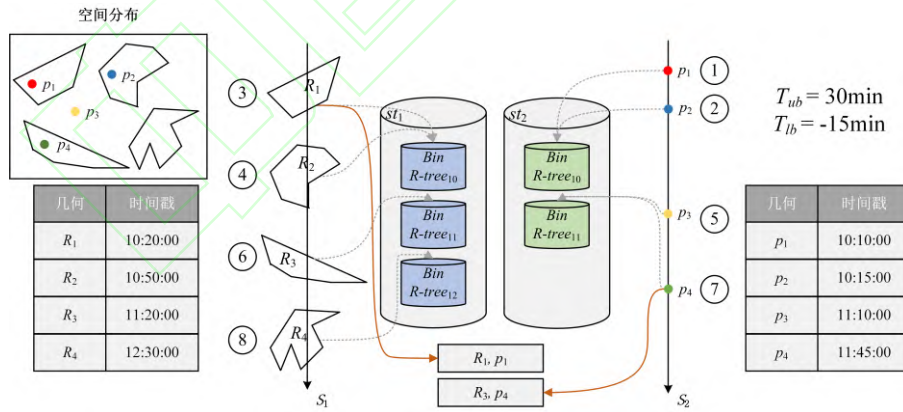


图 10 “流-流”空间间隔连接示例

Fig. 10 Stream-Stream Spatial Interval Join Example

图 10 是一个基于 BSINLJ 算法的示例, 在该示例中  $S_1$  由多边形组成,  $S_2$  由空间点组成, 其中每个几何对象对应的时间戳如图中表格内所示。连接的上界时间为 30 分钟, 下界时间为 15 分钟, 分箱间隔为 1 小时。连接条件为多边形包含空间点。下面基于该示例进一步阐述算法细节:

(1) 元素行为: 元素到达后首先检索对方流中状态

是否有可连接对象, 如果有则连接, 然后将其插入本流的状态中。

(2) 时间间隔基准: 查询的时间间隔以  $S_1$  中的元素为基准,  $S_2$  中的时间间隔需要对称变换为  $[-T_{ub}, -T_{lb}]$ 。示例中,  $S_1$  中 10:20:00 到达的  $R_1$  需要在  $st_2$  中检索  $[10:05:00, 10:50:00]$  的数据, 而  $S_2$  中 10:10:00 到达的  $p_1$  需要在  $st_1$  检索的时间范围是  $[9:40:00, 10:25:00]$ 。



(3) 分箱语义: 本文状态管理基于分箱语义,  $S_1$  和  $S_2$  的状态可能存放于多个 BinR-tree 中, BinR-tree 的建立与删除也以数据驱动。以  $p_3$  为例, 其到达后, 首先需要检索  $st_1$  中 [10:40:00,11:25:00] 的数据, 这些数据分布在  $BinR-tree_{10}$  和  $BinR-tree_{11}$  中, 分别检索; 检索与连接完成后, 需要将  $p_3$  记录  $st_2$ , 其目标  $BinR-tree_{11}$  不存在, 建立并注册其清理时间为 12:15:00。

### 3 实验结果与分析

#### 3.1 实验设置

本文所使用的物理机配备 Intel(R) Xeon(R) Gold 6240 处理器, 其基本频率 2.60 GHz, 最大睿频频率可达 3.90 GHz, 内部含 36 个逻辑内核。内存总量为 256GB, 速率为 2933MT/s。为了模拟分布式环境, 本文在物理机上创建了三个虚拟机用于实验。每个虚拟机分配 8 个逻辑 CPU 内核, 64GB 内存, 使用 Kafka 模拟空间数据流, 虚拟机集群所配备的主要组件及对应版本分别为 JDK 1.8.0\_261、Apache Flink 1.13.3、Apache Kafka 2.6.0。

本文采用了两组不同规模的数据进行实验评估。1) 微软 T-Drive<sup>[26]</sup>数据集。该数据集是由 10,357 辆出租车从 2008 年 2 月 2 日至 8 日在北京市范围内产生的轨迹, 共 17 662 984 条记录; 2) 深圳轨迹数据集。该数据集使用深圳市某年连续 7 天的出租车轨迹数据, 与 T-Drive 数据集同时间跨度, 共 425 205 439 条记录。

本文采用吞吐量和延迟来评价连接查询的性能。吞吐量是指单位时间内系统所处理的数据量, 延迟是指数据记录输入和输出流处理引擎的时间差, 本文对上述指标分别独立测试以保证准确性。在进行吞吐量测试时, 本文在输入端 Kafka 中积累一定数据之后再启动作业执行连接查询, 以避免由于输入端数据不足而导致的 Flink 作业闲置; 在延迟测试时, 则首先启动空间流查询作业, 模拟输入流, 小批量写入数据以减少延迟测试误差。本文针对空间数据流的特性做了一系列查询优化, 为了验证其有效性, 以直接使用 Flink 原生接口的流查询算法 (下文称 Flink 原生方法) 作为基线方法。具体来说, 由于 Flink 未提供空间数据分区与索引支持, 对于“流-表”连接, 遍历维表数据实现连接; 对于“流-流”连接, 采用随机分区 (为保证连接结果不遗漏, 将其中一条流广

播处理), 各分区内遍历所有候选数据进行连接。

#### 3.2 性能评估与分析

##### (1) “流-表”空间连接性能评估

本文将数据集所在地区的行政区划地理多边形作为查询几何, 标记数据流中点的所属行政区。

由图 11 可以看出, 本文所实现的“流-表”空间连接在单并行度下可以达到每秒 6 万以上的吞吐量, 总体延迟平均在 90 毫秒左右, 同原生方法相比, 吞吐量平均有 9 倍提升, 延迟平均降低 20 毫秒, 对于大规模数据查询性能提升更明显。

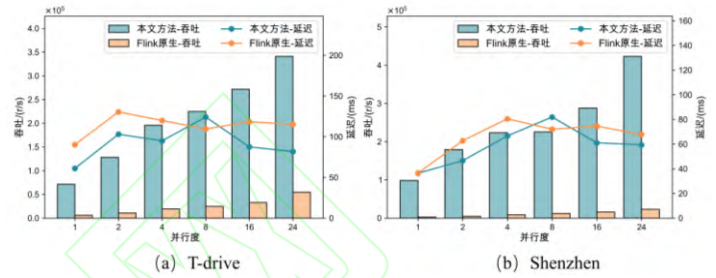


图 11 “流-表”空间连接实验结果

Fig. 11 Experimental Results of Stream-Table Spatial Join

##### (2) “流-流”空间连接性能评估

“流-流”连接实验将数据集模拟两次, 并视其为两个不同的空间数据流, 连接查询条件是距离小于 1 千米。性能将从并行度、窗口 (时间间隔) 大小、网格划分进行实验, 对于间隔连接, 额外进行分箱大小的实验。

##### a. 并行度实验

如图 12 所示, 与“流-表”连接相同, 本文“流-流”连接方法随着并行度增加性能平稳提升, 具有良好的横向拓展性。本文方法查询效率较原生方法有所提升, 且间隔连接的性能提升比窗口连接明显, 这是因为间隔连接不存在窗口语义下流数据的等待行为。

##### b. 窗口 (时间间隔) 实验

从实验结果 (图 13) 来看, 窗口 (时间间隔) 大小对流-流连接的吞吐量影响不大, 吞吐量呈现上升趋势主要是由于查询时间范围扩大导致参与计算的元素增多。而延迟方面, 窗口连接需要根据事件时间的推进, 在特定时间点触发计算, 随着窗口的增大, 连接延迟提升, 远高于流间隔连接。该部分实验也说明了脱离了有界时间语义的流间隔连接更加适合流处理逻辑。

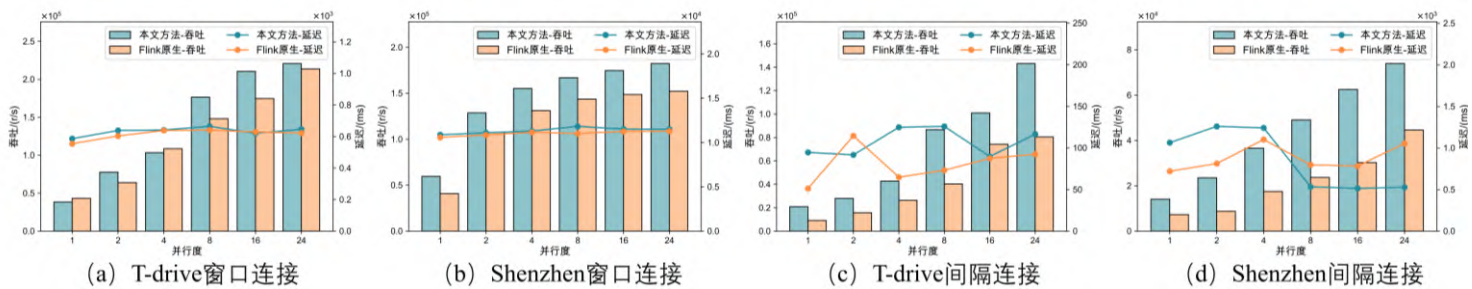


图 12 并行度实验结果

Fig. 12 Parallelism Experiment Results

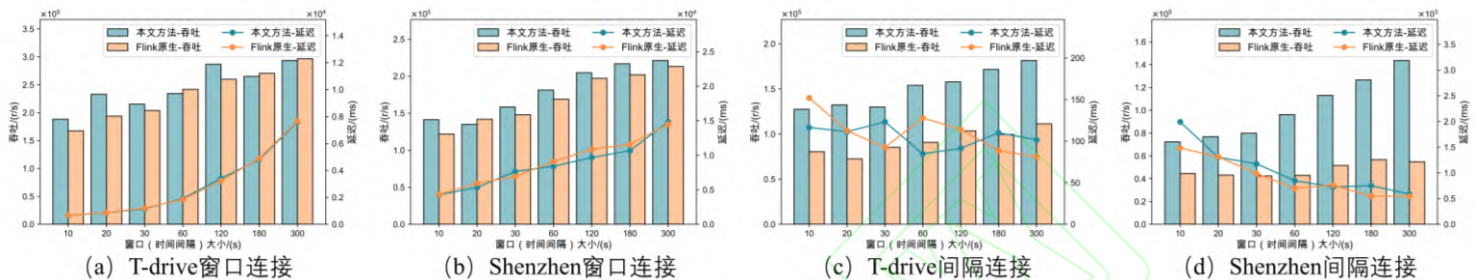


图 13 窗口(时间间隔)实验结果

Fig. 13 Window (Interval) Experiment Results

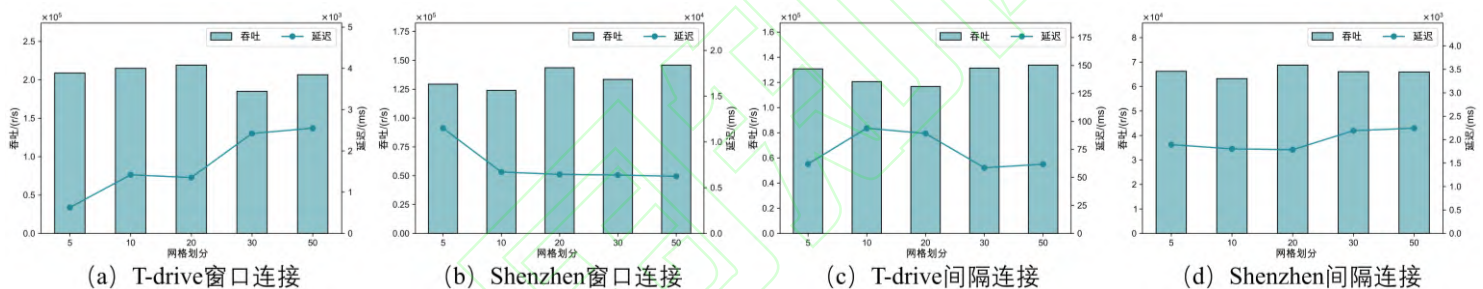


图 14 网格划分实验结果

Fig. 14 Grid Split Experiment Results

### c. 网格划分实验

由于 Flink 原生方法无网格剖分实现, 本部分实验仅进行数据规模对照。实验结果如图 14, 其中, 网格划分  $n$  是指将实验区域剖分为  $n \times n$  个网格, 其主要作用是均衡负载和保证数据的空间邻近性。结果显示, 大部分实验组调整网格划分对于性能影响不大。但是对于 T-drive 数据集, 当网格划分较多后, 通信时间开销和边界冗余的物理开销均会增大, 反而会增加其延迟(见图 14 (a)), 说明网格划分应当根据集群配置和数据分布确定, 不是越多越好。

### d. 分箱大小实验

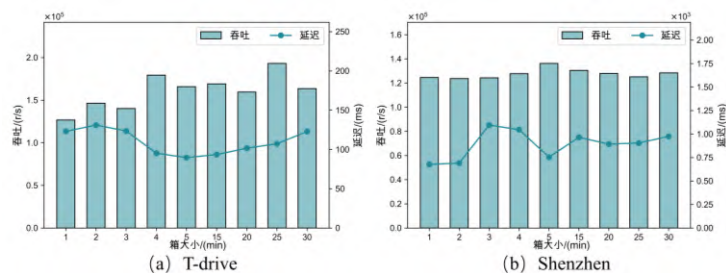


图 15 分箱大小实验结果

Fig. 15 Bin Size Experiment Results

分箱实验统一在 5 分钟时间间隔大小下进行, 实验结果见图 15。对于 T-drive 数据集(图 15 (a)), 当箱大小小于窗口大小时, 查询会跨多个 BinR-tree, 影响查询效率, 此时增大箱大小可以提升查询效率; 当箱大小大于等于时间间隔后, 对查询效率的提升不大。对于 Shenzhen 数据集, 由于其单位时间内数据量更大, 各 BinR-tree 的建立开销小于其查询效率提升缩减的开销, 在实验参数范围内性能变化不大。同时, 在实验过程中也发现, 若箱大小设置过大, 会带来大量过期数据的堆



积,甚至因为内存占用过大导致流作业执行失败,将状态存储至外部系统是提高容错性的方法。

## 4 结 语

本文聚焦空间数据流的连接查询问题,基于对连接场景共性问题的分析,形式化定义了针对空间数据流各类连接查询处理问题,包括针对空间数据流和空间维度表的“流-表”空间连接,以及针对两条空间数据流的“流-流”空间窗口连接和“流-流”空间间隔连接。在此基础上,基于全局网格数据分区及局部空间索引的两层数据索引结构,针对不同连接场景分别设计了适应的连接策略:对于“流-表”连接,提出了基于两级 R-tree 的拓扑关系判断优化算法;对于“流-流”连接,本文设计了一种分区边界数据冗余算法,以保证分区边界数据的正确连接。此外,对于间隔时间语义的缓存需求,基于数据分箱思想提出了一种兼顾状态管理与检索效率的 BinR-tree 结构。最后,在 Apache Flink 中实现了上述连接查询算法,并基于模拟数据流开展了大量性能测试。实验结果表明,本文所定义的连接处理适用于实时空间数据流处理场景,相较基线方法具有不同程度的性能提升。

在后续工作中,将对空间数据流查询的优化方法进一步研究,包括设计更低维护成本的空间流索引结构、设计更多细化场景的空间流连接算子(如多流连接算子)、对于更多拓扑判断规则的优化、以及提高大规模空间数据流查询处理的容错性。

## 参 考 文 献

[1] Sun Dawei, Zhang Guangyan, Zheng Weimin. Big Data Streaming Computing: Technologies and instances[J]. Ruan Jian Xue Bao/Journal of Software, 2014, 25(4): 839-862 (孙大为, 张广艳, 郑纬民. 大数据流式计算: 关键技术及系统实例[J]. 软件学报, 2014, 25(4): 839-862)

[2] ZHU Jianjun, SONG Yingchun, HU Jun, ZOU Bin, WU Lixin. Challenges and Development of Data Processing Theory in the Era of Surveying and Mapping Big Data[J]. Geomatics and Information Science of Wuhan University, 2021, 46(7): 1025-1031. doi: 10.13203/j.whugis20210232(朱建军, 宋迎春, 胡俊, 等. 测绘大数据时代数据处理理论面临的挑战与发展[J]. 武汉大学学报(信息科学版), 2021, 46(7): 1025-1031.)

[3] LI Jun, LIU Juqing, ZHAO Xuesheng, HUANG Qian, SUN Wenbin, XU Zhihua, WANG Hao. Trajectory

Data Management and Analysis Framework Based on Geographical Grid Model: Method and Application[J]. Geomatics and Information Science of Wuhan University, 2021, 46(5): 640-649. doi: 10.13203/j.whugis20200459(李军, 刘举庆, 赵学胜, 等. 地理网格模型支持下的轨迹数据管理与分析框架: 方法与应用[J]. 武汉大学学报(信息科学版), 2021, 46(5): 640-649.)

[4] CHEN Dong, ZHANG Xiang, CHEN Nengcheng. Smart City Awareness Base Station: A Prospective Integrated Sensing Infrastructure for Future Cities[J]. Geomatics and Information Science of Wuhan University, 2022, 47(2): 159-180. doi: 10.13203/j.whugis20210224(陈栋, 张翔, 陈能成. 智慧城市感知基站: 未来智慧城市的综合感知基础设施[J]. 武汉大学学报(信息科学版), 2022, 47(2): 159-180.)

[5] QIU A'gen. In-memory Distributed Computing Based Approximate Query Processing on Spatial Data[J]. Acta Geodaetica et Cartographica Sinica, 2017, 46(12): 2044-2044.(仇阿根. 基于分布式内存计算的空间数据近似查询处理方法[J]. 测绘学报, 2017, 46(12): 2044-2044.)

[6] SHENG Yuyu, BI Shuoben, FAN Jingjin, NKUNZIMANA Athanase, XU Zhihui. Analyzing Spatiotemporal Patterns of Traffic Hotspots Using Traffic Operation Indicators[J]. Geomatics and Information Science of Wuhan University, 2021, 46(5): 746-754. doi: 10.13203/j.whugis20190357(盛宇裕, 毕硕本, 范京津, 等. 运用交通运行状况指标分析交通热点时空模式[J]. 武汉大学学报(信息科学版), 2021, 46(5): 746-754.)

[7] Shapiro L D. Join processing in database systems with large main memories[J]. ACM Transactions on Database Systems (TODS), 1986, 11(3): 239-264

[8] Aji A, Wang F, Vo H, et al. Hadoop gis: a high performance spatial data warehousing system over mapreduce[J]. Proceedings of the VLDB Endowment, 2013, 6(11): 1009-1020.

[9] Eldawy A, Mokbel M F. Spatialhadoop: A mapreduce framework for spatial data[C]//2015 IEEE 31st international conference on Data Engineering. IEEE, 2015: 1352-1363

[10] You S, Zhang J, Gruenwald L. Large-scale spatial join query processing in cloud[C]//2015 31st IEEE international conference on data engineering



- workshops. IEEE, 2015: 34-41.
- [11] Yu J, Wu J, Sarwat M. Geospark: A cluster computing framework for processing large-scale spatial data[C]//Proceedings of the 23rd SIGSPATIAL international conference on advances in geographic information systems. 2015: 1-4.
- [12] Yu J, Zhang Z, Sarwat M. Spatial data management in apache spark: the geospark perspective and beyond[J]. *GeoInformatica*, 2019, 23(1): 37-78.
- [13] Xie D, Li F, Yao B, et al. Simba: Efficient in-memory spatial analytics[C]//Proceedings of the 2016 international conference on management of data. 2016: 1071-1085.
- [14] Tang M, Yu Y, Malluhi Q M, et al. Locationspark: A distributed in-memory data management system for big spatial data[J]. *Proceedings of the VLDB Endowment*, 2016, 9(13): 1565-1568.
- [15] Pandey V, Kipf A, Neumann T, et al. How good are modern spatial analytics systems?[J]. *Proceedings of the VLDB Endowment*, 2018, 11(11): 1661-1673
- [16] Shahvarani A, Jacobsen H A. Distributed stream KNN join[C]//Proceedings of the 2021 International Conference on Management of Data. 2021: 1597-1609.
- [17] Zhang F, Zheng Y, Xu D, et al. Real-time spatial queries for moving objects using storm topology[J]. *ISPRS International Journal of Geo-Information*, 2016, 5(10): 178.
- [18] Galić Z, Mešković E, Osmanović D. Distributed processing of big mobility data as spatio-temporal data streams[J]. *Geoinformatica*, 2017, 21: 263-291.
- [19] Shaikh S A, Mariam K, Kitagawa H, et al. GeoFlink: A distributed and scalable framework for the real-time processing of spatial streams[C]//Proceedings of the 29th ACM International Conference on Information & Knowledge Management. 2020: 3149-3156.
- [20] Shaikh S A, Kitagawa H, Matono A, et al. GeoFlink: An Efficient and Scalable Spatial Data Stream Management System[J]. *IEEE Access*, 2022, 10: 24909-24935.
- [21] *Encyclopedia of GIS*[M]. Springer Science & Business Media, 2007.
- [22] Ross M, Kimball R. *The data warehouse toolkit: the definitive guide to dimensional modeling*[M]. John Wiley & Sons, 2013.
- [23] Tian R, Zhai H, Zhang W, et al. A Survey of Spatio-Temporal Big Data Indexing Methods in Distributed Environment[J]. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2022, 15: 4132-4155.
- [24] Hughes J N, Annex A, Eichelberger C N, et al. Geomesa: a distributed architecture for spatio-temporal fusion[C]//Geospatial informatics, fusion, and motion video analytics V. SPIE, 2015, 9473: 128-140.
- [25] Foley J D, Van F D, Van Dam A, et al. *Computer graphics: principles and practice*[M]. Addison-Wesley Professional, 1996.
- [26] Yuan J, Zheng Y, Zhang C, et al. T-drive: driving directions based on taxi trajectories[C]//Proceedings of the 18th SIGSPATIAL International conference on advances in geographic information systems. 2010: 99-108.

#### 网络首发:

**标题:** 面向大规模空间数据流的分布式连接查询方法

**作者:** 李柞霖, 向隆刚, 余列冰, 吴华意, 关雪峰

**Doi:** 10.13203/j.whugis20230040

**收稿日期:** 2023-09-11

#### 引用格式:

李柞霖, 向隆刚, 余列冰, 等. 面向大规模空间数据流的分布式连接查询方法[J]. *武汉大学学报 (信息科学版)*, 2023, Doi:10.13203/j.whugis20230040. (LI Zuolin, XIANG Longgang, YU Liebing, et al. Distributed Join Query Method for Large-Scale Spatial Data Streams[J] *Geomatics and Information Science of Wuhan University*, 2023, Doi:10.13203/j.whugis20230040.)