

RFID Data Processing in Supply Chain Management Using a Path Encoding Scheme

Chun-Hee Lee and Chin-Wan Chung

Abstract—RFID technology can be applied to a broad range of areas. In particular, RFID is very useful in the area of business, such as supply chain management. However, the amount of RFID data in such an environment is huge. Therefore, much time is needed to extract valuable information from RFID data for supply chain management. In this paper, we present an efficient method to process a massive amount of RFID data for supply chain management. We first define query templates to analyze the supply chain. We then propose an effective path encoding scheme that encodes the flows of products. However, if the flows are long, the numbers in the path encoding scheme that correspond to the flows will be very large. We solve this by providing a method that divides flows. To retrieve the time information for products efficiently, we utilize a numbering scheme for the XML area. Based on the path encoding scheme and the numbering scheme, we devise a storage scheme that can process tracking queries and path oriented queries efficiently on an RDBMS. Finally, we propose a method that translates the queries to SQL queries. Experimental results show that our approach can process the queries efficiently.

Index Terms—RFID, supply chain management, path encoding scheme, prime number.

1 INTRODUCTION

As the size of the RFID tag becomes smaller and the price of the RFID tag gets lower, RFID technology has been applied to many areas. It is different from existing technologies, such as barcode systems and magnetic card systems, in that barcode systems and magnetic card systems require contact between a detector and an object. RFID readers in RFID systems, on the other hand, can detect RFID tags without contact.

A typical example of RFID technology being used is supply chain management. In supply chain management, in order to know the movements of products easily, an RFID tag is attached to a product. If the product with an RFID tag moves or stays near the detection region, RFID readers will detect RFID tags and the detected information will be generated in the form of (*tag identifier*,¹ *location*, *time*). As the flow of the product is detected easily by RFID technology, it is observed that RFID can be used to revolutionize supply chain management.

The RFID data generated in each region (i.e., (*tag identifier*, *location*, *time*)) are sent to the central server. Then, the data are transformed into stay records in the form of (*tag identifier*, *location*, *start time*, *end time*). While raw RFID data have many duplicates, the transformed data (i.e., stay records) do not have duplicates. We can represent how long a tag stays at a

location by the *start time* and *end time* of stay records. The stay records for each tag compose a trace record that gives us movement history with time information for the tag. In this paper, we will use trace records instead of stay records for storing RFID data in the central server.

We can store RFID data in a relational table BASIC_TABLE(*TAG_ID*, *LOC*, *START_TIME*, *END_TIME*) as a straightforward method, where *TAG_ID* represents the tag identifier, *LOC* the location, *START_TIME* the time when the tag enters the location, and *END_TIME* the time when the tag leaves the location. The queries that analyze the supply chain are related to the product transition. For example, a manager may ask the query “Find the number of laptops that go through locations *Factory1*, *Distribution Center1*, and *Store1*.” To evaluate the query with BASIC_TABLE, we must perform self-joins of BASIC_TABLE many times. Also, the size of BASIC_TABLE is big. Therefore, it requires a lot of time to execute the query with BASIC_TABLE.

To support efficient path-dependent aggregates for RFID data, Gonzalez et al. propose a new warehousing model [14]. In the model, STAY_TABLE(*GID*, *LOC*, *START_TIME*, *END_TIME*, *COUNT*) is used to store RFID data efficiently. In many RFID applications, products usually move together in large groups during the early stages and move in small groups during the later stages. Therefore, they represent stay records with the same location and time by one record such as (*tag identifier list*, *location*, *start time*, *end time*, *the number of tags with the same location and time*). To link locations efficiently (i.e., to perform self-joins of STAY_TABLE efficiently), the tag identifier list is encoded by the prefix encoding scheme. The value of the prefix encoding scheme corresponds to *GID* in STAY_TABLE. To know whether two locations (A and B) are linked, the prefix encoding scheme checks whether the *GID* that corresponds to A is the prefix of the *GID* that corresponds to B.

Therefore, the approach by Gonzalez et al. [14] reduces the size of the table significantly and improves the join cost. However, if products do not move together in large groups,

1. As a tag identifier, RFID systems use EPC (Electronic Product Code) [2], which is a coding scheme of RFID tags, to identify the tags uniquely.

• The authors are with the Department of Computer Science, Korea Advanced Institute of Science and Technology (KAIST), 291 Daehak-ro, Yuseong-gu, Daejeon 305-701, Republic of Korea.
E-mail: leechun@islab.kaist.ac.kr, chungcw@kaist.edu.

Manuscript received 3 Nov. 2008; revised 10 July 2009; accepted 14 Nov. 2009; published online 12 Aug. 2010.

Recommended for acceptance by Q. Li.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number TKDE-2008-11-0575. Digital Object Identifier no. 10.1109/TKDE.2010.136.

the size of the table will not be largely reduced. In this case, the performance of the approach by Gonzalez et al. does not have a great benefit compared to that of the approach with BASIC_TABLE. Also, since the prefix encoding scheme uses a string comparison in joining tables, it needs much more time than a number comparison. Thus, we propose a new approach to store RFID data and process queries for supply chain management. Since the queries related to the object transition in supply chain management was not defined formally in [14], we first define query templates for tracking queries and path-oriented queries to analyze the supply chain. Then, to solve the above drawbacks, we propose a new approach.

While Gonzalez et al. [14] focus on groups in which products move together, we focus on the movement for each tag. The movement of one tag makes a path in supply chain management. Therefore, we devise a path encoding scheme to process tracking queries and path-oriented queries efficiently. The proposed path encoding scheme can encode a path with only two numbers, which is motivated by Wu et al. [39]. The numbers are called Element List Encoding Number and Order Encoding Number, and will be explained in detail subsequently. In [39], in order to determine whether a relationship exists between two elements in an XML tree, a property of prime numbers is used. In addition, to preserve the global order for elements, simultaneous congruence values are used. In this paper, we use the property of prime numbers to encode nodes in a path, and simultaneous congruence values to encode the ordering among nodes in the path. Our encoding scheme is based on the Fundamental Theorem of Arithmetic and the Chinese Remainder Theorem. Using the proposed path encoding scheme, we can efficiently retrieve paths that satisfy the path condition in a query. To store the time information related to the movement, we separate the time information from trace records. We use the region numbering scheme [40], which is widely used in the XML area, in order to retrieve the time information efficiently. Based on the path encoding scheme and the region numbering scheme, we devise a new relational schema to store the path information and the time information for tags.

If a path is long, its Element List Encoding Number and Order Encoding Number are very large and an RDBMS may not support the data type for storing very large numbers. Although we can implement our encoding scheme on a special purpose database engine, it has many disadvantages. Therefore, we propose a method to divide a path. If we cannot store a path by using one attribute in an RDBMS, we divide the path into multiple path segments, each of which can be stored by using one attribute.

Based on the encoding schemes described above, we translate query templates into SQL queries. However, in the case of dividing paths, the SQL translation is not straightforward due to the computation overflow. Therefore, we propose a method to translate query templates into SQL queries, which do not cause the computation overflow, by using mathematical analyses.

1.1 Contributions

The contributions of the paper are as follows:

- **Encoding Scheme to Encode Flows of Products:** To support tracking queries and path-oriented queries

efficiently, we propose a path encoding scheme for flows of products. Element List Encoding Number is computed by the product of the prime numbers that correspond to the nodes in a path. Based on the Chinese Remainder Theorem, Order Encoding Number is computed. By using the two numbers, we can easily find paths that satisfy the path condition.

- **Efficient Relational Schema and Query Translation:** We propose an efficient relational schema with the path encoding scheme and the region numbering scheme. Based on the schema, we propose a method that translates tracking queries and path-oriented queries into SQL queries.
- **Effective Method to Store Long Paths:** A long path is common in typical supply chain management. The product of prime numbers for a long path can become excessively large, which can cause a limit to the path length. To avoid overflows in storing large numbers, we provide a method that divides the long path into several path segments.
- **Extended Query Translation for Handling Long Paths:** When we store a long path by dividing the path, the computation overflow may occur during the execution of the translated SQL query. To overcome the computation overflow, we extend the method to translate query templates into SQL queries by using mathematical formulas.
- **Defining of Query Templates to Analyze the Supply Chain:** We define query templates to analyze the supply chain. We consider query templates for tracking queries and path-oriented queries. For path-oriented queries, we provide a grammar to effectively express the path condition for products like XPath [3].
- **Experimentation to Validate Our Proposed Approach:** Through extensive experiments, we show that our approach is efficient. Experimental results show that the query performance of our approach is better than the recent approach in most queries.

1.2 Organization

The rest of the paper is organized as follows: In Section 2, we discuss related work on managing RFID data. In Section 3, we deal with data formats and define query templates for tracking queries and path-oriented queries in supply chain management. We show the architecture that stores RFID data and process queries in Section 4. We describe our path encoding scheme in Section 5 and devise a storage scheme based on our path encoding scheme in Section 6. We propose a method to translate tracking queries and path-oriented queries into SQL queries in Section 7, and experimentally show the superiority of our approach in Section 8. We make a conclusion in Section 9.

2 RELATED WORK

In contrast to the initial study for RFID, which focuses on the device, many studies have been done recently to manage RFID data as the amount of RFID data has become large.

The system architecture for managing RFID data is discussed in [8], [11], [17]. Bornhövd et al. [8] describe the

Auto-ID infrastructure (Device Layer/Device Operation Layer/Business Processing Bridging Layer/Enterprise Application Layer). The Auto-ID infrastructure integrates data from smart items such as RFID and sensor with existing business processes. Chawathe et al. [11] suggest a layered architecture for managing RFID data (Tag/Reader/Savant/EPC-IS/ONS server). In the work of Hag et al. [17], two software layers are proposed: the ubiquity agent architecture and the tag centric RFID application architecture. The ubiquity agent architecture is for a general-purpose core architecture. The tag centric RFID application architecture is for the RFID application architecture that specializes the generic agent architecture.

RFID data are generated in the form of streaming data and then stored in a database for data analyses. Therefore, there are two types of approaches for managing RFID data. One approach is online processing for RFID data and it is related to data stream processing [5], [6], [19], [36]. The other approach is offline processing and it is related to stored data processing [4], [7], [13], [14], [18], [29], [35].

In the aspect of viewing RFID data as data streams, event processing and data cleaning have been studied. RFID data have a temporal property, which is important in analyzing data. Therefore, a temporal RFID event can be defined. However, this cannot be well supported by traditional ECA (Event-Condition-Action) rule systems. Therefore, Wang et al. [36] formalize the specification and semantics of RFID events and rules including temporal RFID events. Also, they propose a method to detect RFID complex events efficiently. Bai et al. [6] explore the limitation of SQL queries in supporting temporal event detection and discuss an SQL-based stream query language to provide comprehensive temporal event detection.

Inevitably, RFID data have some errors, such as duplicate readings and missing readings. To rescue missing readings, the first declarative and adaptive smoothing filter for RFID data (SMURF) [19] is proposed. SMURF controls the window size of the smoothing filter adaptively using statistical sampling. Also, Bai et al. [5] propose several methods to filter RFID data.

In the aspect of viewing RFID data as stored data, various approaches exist that can manage RFID data. Since RFID readers can detect RFID tags easily and quickly, object tracking that uses RFID is widely used. To trace tag locations, a new index is proposed by Ban et al. [7]. They point out the problem of representing trajectories in RFID data and propose a new data model to solve it. Also, they devise a new index scheme called the Time Parameterized Interval R-Tree as a variant of the R-Tree. In order to represent a collection of tag identifiers generated by item tracking applications compactly, a bitmap data type is proposed and a set of bitmap access and manipulation routines is provided in [18]. Agrawal et al. [4] deal with the tracing of items in distributed RFID databases. They introduce the concept of traceability networks and propose an architecture for traceability query processing in distributed RFID databases.

Since RFID data have a temporal property, it is difficult to model RFID data by using the traditional ER-model. Therefore, Wang et al. [35] propose a new model called the

Dynamic Relationship ER Model (DRER), which simply adds a new relationship (dynamic relationship). They also propose methods to express temporal queries based on DRER by using SQL queries. Although we can use the above techniques, such as indexes and bitmap data types, to process tracking queries and path-oriented queries, it is inefficient to process the queries since they do not consider the object transition.

Gonzalez et al. [14] propose a new warehousing model for the object transition and a method to process a path selection query. To get aggregate measures on a path, they join tables many times. Therefore, they use compression in order to reduce the join cost. However, the proposed compression is useless if products do not move together in large groups.

A preliminary version of this work appeared in [22]. While an efficient storage scheme and query processing for RFID data in supply chain management were proposed by Lee and Chung [22], they cannot be applied when the path length is large since Element List Encoding Number and Order Encoding Number cannot be stored by using the data types that an RDBMS supports. Long paths are common in typical supply chain management because there are many checkpoints in a route. In such a case, the work of Lee and Chung [22] should be implemented on a special-purpose DBMS that can support very large numbers.

To solve the problem, we propose a method to divide a long path. By dividing the path, we can store the long path without modifying an RDBMS or implementing a special-purpose database engine. However, in the case of dividing the path, if we use the translation method proposed in [22], computation overflow may occur. Therefore, we provide a method to translate query templates into SQL queries, which do not cause the computation overflow, by mathematical analyses. Also, we redo the entire experiment with data sets generated from a real-life example and additional queries.

In addition, managing object transitions has been dealt with in various areas with different views and terminologies. In the vehicle tracking area, most papers focus on how to detect and track vehicle objects than how to manage the generated tracking data. To track vehicles, video cameras are generally used [16], [20], [21]. To detect and track vehicles from video streams efficiently, Kanhere et al. [20] and Kim et al. [21] propose image-processing-based methods while Haidarian-Shahri et al. [16] use a graph-based approach. They do not consider how to manage the tracking data for vehicles. In this paper, we focus on managing the flows of products that are similar to the tracking data.

In XML databases, the path information can be stored in XML data format. In [10], [25], [34], XML data can be represented as compact data (i.e., compressed data) and direct query processing is available in such a compact representation. Also, to index and query XML data, various approaches are proposed [23], [24], [30], [33], [37], [39], [40]. However, in an environment where so much path data are generated, the existing approaches in the XML area cannot manage and analyze the path information for products efficiently since they do not focus on many ancestor/descendant relationships in a query.

In spatio-temporal databases, the data generated from moving objects (e.g., vehicles) are manipulated. As moving

(tag1, A, 2), (tag4, A, 2), (tag2, A, 2), (tag3, A, 2), (tag1, A, 3),	tag1: A[2,3]->B[5,7]->C[8,9]
(tag2, A, 3), (tag4, A, 3), (tag3, A, 3), (tag3, B, 5), (tag1, B, 5),	tag2: A[2,3]->B[5,7]->C[8,9]
(tag2, B, 5), (tag4, B, 5), (tag1, B, 6), (tag4, B, 6), (tag3, B, 7),	tag3: A[2,3]->B[5,7]->C[8,9]
(tag1, B, 7), (tag2, B, 7), (tag4, B, 7), (tag2, C, 8), (tag1, C, 8),	tag4: A[2,3]->B[5,7]->D[13,16]
(tag3, C, 8), (tag3, C, 9), (tag1, C, 9), (tag2, C, 9), (tag4, C, 13),	
(tag4, C, 14), (tag4, C, 16)	
(a)	(b)

Fig. 1. (a) Raw data and (b) trace records.

objects go around, the object movement information is generated and its collection composes the trajectories of moving objects. In [12], [15], the data structure and model for moving objects are proposed. Also, in [26], [27], [28], [32], [38], to index and query the data from moving objects, various index structures and algorithms are used. In the aspect of representing object movements, the trajectory data of moving objects and the flow information of products in supply chain management using RFID are similar. Both the trajectories of moving objects and the flow information of products can be represented as the set of $\langle \text{ObjectID}, \text{Location}, \text{Time} \rangle$. However, the locations of products in RFID applications are limited since they are detected through the fixed RFID readers. In contrast, the locations of moving objects in spatio-temporal databases are not limited. Also, in spatio-temporal databases, most queries are related to the locations of objects (e.g., k-nearest neighbor query, range query) rather than the flow of objects. Therefore, query processings in these two areas are quite different, although the data can be represented in a similar format.

3 PROBLEM DEFINITION

Raw RFID data consist of a set of triples $(\text{TagID}, \text{Loc}, \text{Time})$, where

- TagID is the Electronic Product Code (EPC) of the tag and is used for identifying the tag uniquely,
- Loc is the location of the RFID reader which detects the tag,
- Time is the time of detecting the tag.

We translate raw RFID data generated in supply chain management into a set of stay records that do not have duplicates. A stay record has the form $(\text{TagID}, \text{Loc}, \text{StartTime}, \text{EndTime})$, where

- TagID and Loc are the same as above,
- StartTime is the time when the tag enters the location,
- EndTime is the time when the tag leaves the location.

From the stay records of a tag, we can construct the trace record of the tag in the form of $\text{TagID}: L_1[S_1, E_1] \rightarrow \dots \rightarrow L_n[S_n, E_n]$, where L_1, \dots, L_n are the locations where the tag is detected, S_i is StartTime at the location L_i , E_i is EndTime at the location L_i , and $L_i[S_i, E_i]$ is ordered by S_i . We use a set of trace records instead of raw RFID data in our systems.

Example 1. Fig. 1a shows raw RFID data and Fig. 1b shows a set of trace records that corresponds to raw

RFID data. From (tag1, A, 2) and (tag1, A, 3), we get the stay record (tag1, A, 2, 3). Similarly, we can compute stay records (tag1, B, 5, 7) and (tag1, C, 8, 9) for tag1. Finally, we can compute the trace record, $\text{tag1} : A[2, 3] \rightarrow B[5, 7] \rightarrow C[8, 9]$.

To analyze the supply chain, we use queries for the object transition. Although Gonzalez et al. [14] use a path selection query, it is insufficient to express the relationship between locations. Therefore, we define query templates to analyze the supply chain. We consider query templates for tracking queries and path-oriented queries. The tracking query finds the movement history for the given tag. The path-oriented query is classified into the path-oriented retrieval query and the path-oriented aggregate query. The path-oriented retrieval query finds tags that satisfy given conditions (including a path condition), and the path-oriented aggregate query computes the aggregate value for tags that satisfy given conditions (including a path condition). In query templates for the path-oriented retrieval query and the path-oriented aggregate query, we provide a grammar to effectively express path conditions for products like XPath [3].

Fig. 2 shows the formal definition for query templates in supply chain management. There are three query templates (tracking query, path-oriented retrieval query, path-oriented aggregate query). The query template for a tracking query has only a tag identifier to trace the tag. The path-oriented retrieval query consists of Path Condition and Info Condition. The path-oriented aggregate query needs Aggregate Function as well as Path Condition and Info Condition. Path Condition uses a grammar similar to XPath. Path Condition consists of a Step sequence. Each Step has a parent axis (/) or an ancestor axis (/ /). Also, each Step may have Time Conditions. Time Condition is the predicate for StartTime and EndTime . The argument for Aggregate Function (i.e., Time Selection) allows only the time information. We can express various queries in supply chain management by using the query templates in Fig. 2. We show some examples for tracking queries and path-oriented queries in Fig. 3.

Our problem definition is as follows: For an environment where there is a large amount of RFID data in supply chain management and users issue tracking queries and path-oriented queries, devise an efficient storage scheme and processing method for the queries.

4 ARCHITECTURE

Fig. 4 shows the architecture to store RFID data, and process tracking queries and path-oriented queries in supply chain

tag1: A[2,3]->B[5,7]->C[8,9]	tag6: A[2,3]->E[4,6]->C[7,8]
tag2: A[2,3]->B[5,7]->C[8,9]	tag7: A[2,3]->E[4,6]->C[7,8]
tag3: A[2,3]->B[5,7]->C[8,9]	tag8: A[2,3]->E[4,6]
tag4: A[2,3]->B[5,7]->D[13,16]	tag9: A[2,3]->D[4,5]
tag5: A[2,3]->B[7,8]->D[14,18]	tag10: A[2,3]->D[5,6]

Fig. 5. Example of trace records.

with cycles later. The path of each tag composes the tree by eliminating duplicate nodes. The numbers beside the nodes are prime numbers, which will be explained subsequently. In Fig. 6, the dark nodes mean that there are tags whose final location is the node. We store all paths that end at dark nodes, which are $A \rightarrow B \rightarrow C$, $A \rightarrow B \rightarrow D$, $A \rightarrow E$, $A \rightarrow E \rightarrow C$, and $A \rightarrow D$. Although a huge amount of RFID data is generated, the size of the tree is small.

To encode a path, we can consider various techniques in the XML area [23], [25], [30], [33], [37], [39], [40]. However, those techniques are inefficient in processing queries that have many ancestor-descendant relationships, such as “Find the tags that go through locations L_1, L_2, L_3 ($//L_1//L_2//L_3$).” In supply chain management, we need such queries to analyze the flows of tags. Before proposing a new path encoding scheme, we introduce two well-known theorems [31].

Theorem 1. *The Fundamental Theorem of Arithmetic (The Unique Factorization Theorem): Any natural number greater than 1 is uniquely expressed by the product of prime numbers.*

For example, $231 = 3 \times 7 \times 11$ and the product of any other prime number combination for 231 does not exist.

Theorem 2. *The Chinese Remainder Theorem: Suppose that p_1, p_2, \dots , and p_n are pairwise relatively prime numbers. Then, there exists X with $0 \leq X < p_1 p_2 \dots p_n$ which solves the system of simultaneous congruences.*

$$\begin{aligned} X \bmod p_1 &= a_1 \\ X \bmod p_2 &= a_2 \\ &\dots \\ X \bmod p_n &= a_n \end{aligned}$$

For example, consider the system of simultaneous congruences, such as $X \bmod 3 = 2$, $X \bmod 7 = 3$, and $X \bmod 11 = 2$. Then, by Theorem 2, there exists X with $0 \leq X < 3 \times 7 \times 11$. In this example, X is 101. We can compute X by using the method in [1].

Let $L_1 \rightarrow L_2 \rightarrow \dots \rightarrow L_n$ be a path to encode. Suppose that each location is associated with a different prime number, the nodes with the same location have the same prime number, and the prime number for location L_i is denoted by $Prime(L_i)$. Then, we define Element List Encoding Number for the path $L_1 \rightarrow L_2 \rightarrow \dots \rightarrow L_n$ as $Prime(L_1) \times Prime(L_2) \times \dots \times Prime(L_n)$. If Element List Encoding Number is given, we can know the locations that compose the path since Element List Encoding Number is uniquely factorized by the product of prime numbers that correspond to locations by Theorem 1. However, although we know the locations in the path by Element List Encoding Number, we cannot know the ordering between the locations. To encode the ordering

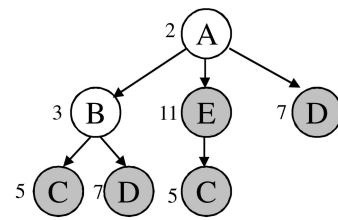


Fig. 6. Tree structure for the trace records.

information compactly and efficiently, we consider the system of simultaneous congruences.

$$\begin{aligned} X \bmod Prime(L_1) &= 1 \\ X \bmod Prime(L_2) &= 2 \\ &\dots \\ X \bmod Prime(L_n) &= n \end{aligned}$$

Here, $1, 2, \dots$, and n are the levels (i.e., orders) of the nodes that correspond to L_1, L_2, \dots , and L_n , respectively. Since $Prime(L_1), Prime(L_2), \dots$, and $Prime(L_n)$ are prime numbers, they are pairwise relatively prime numbers. Thus, by Theorem 2, there exists X with $0 \leq X < Prime(L_1) \times Prime(L_2) \times \dots \times Prime(L_n)$ which solves the system of the above simultaneous congruences. We call X Order Encoding Number. Given Order Encoding Number, we can know the order information for any location L_i in the path by computing $X \bmod Prime(L_i)$. Sometimes, $Prime(L_i)$ may be less than i . To prevent it, a prime number that is greater than the maximum path length is assigned to a location.

Our use of prime numbers is similar to that in [39]. However, the approach in [39] assigns different prime numbers to each element in an XML tree that can have millions of elements, while our approach assigns different prime numbers to different locations in a tree for trace records that has, at most, a few hundred locations. Therefore, in a typical application, the approach in [39] generates extremely large prime numbers. Furthermore, the approach in [39] orders all the elements of an XML tree, while we order only the locations on a path whose length is much smaller than the tree for trace records.

Therefore, we can encode the path by using Element List Encoding Number and Order Encoding Number. Although a path condition has multiple ancestor-descendant relationships, we can find out whether a path satisfies the path condition in Fig. 2 by checking some simple mathematical conditions. We will explain how we process tracking queries and path-oriented queries efficiently by using Element List Encoding Number and Order Encoding Number in Section 7.

Example 2. Assume that in Fig. 6, the prime numbers that correspond to A, B, C, D , and E are 2, 3, 5, 7, and 11. Consider the path $A \rightarrow B \rightarrow C$. Element List Encoding Number for the path is $2 \times 3 \times 5 = 30$. To compute Order Encoding Number, we must compute X such that $X \bmod 2 = 1, X \bmod 3 = 2$, and $X \bmod 5 = 3$. By Theorem 2, there exists X with $0 \leq X < 30$. In this case, $X = 23$. Similarly, we encode the other four paths, and get Element List Encoding Numbers and Order Encoding Numbers for the paths.

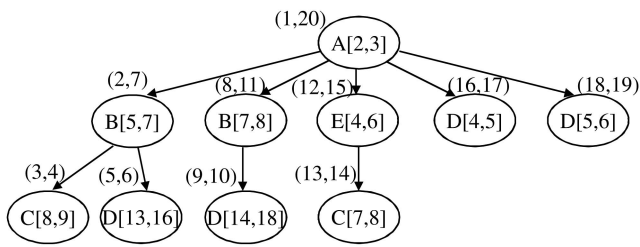


Fig. 7. Time tree structure for the trace records.

Consider a path with cycles. In such a case, we cannot compute Order Encoding Number since the same prime numbers should return different orders. However, we can solve the problem with a simple renaming method that adds the sequence order to the location. For example, consider the path $A \rightarrow B \rightarrow C \rightarrow A \rightarrow B \rightarrow D \rightarrow A$. In the path, the location A appears in the first, fourth, and seventh positions. We can rename A in the first, fourth, and seventh positions to A_1 , A_2 , and A_3 , respectively, by the sequence order of A . In the same way, we can rename the path to $A_1 \rightarrow B_1 \rightarrow C_1 \rightarrow A_2 \rightarrow B_2 \rightarrow D_1 \rightarrow A_3$. Then, since the nodes with the same location name are considered differently, we can compute Order Encoding Number. In data sets that have cycles, a user can write a query with a sequence order such as $//A_1//B_1//A_3$. If a user does not specify the sequence order, we translate the location name with the default sequence order ($= 1$). For example, $//A//D$ is translated into $//A_1//D_1$.

Even though we can store the flow information for products effectively by using the path encoding scheme, we have not discussed the time information for products yet. To store the time information for products, we construct the time tree from trace records in which the node has the start time and end time as well as the location.

In the time tree, we say that two paths are the same if the flows for locations are the same as well as the time information (start time and end time) for locations is the same.

Fig. 7 shows the time tree constructed from the trace records in Fig. 5. The construction of the time tree is the same as that of the tree for the trace records except that the stay records with the same location become different nodes if their time information is different. See $B[5,7]$ node and $B[7,8]$ node in Fig. 7. Although these two nodes have the same location, they are classified as different nodes since they have different time information.

To retrieve the time information efficiently, we store numbers with nodes in the time tree by using the region

numbering scheme [40] which assigns a node two values, Start and End. Start and End are assigned consecutively during the depth-first search. The region numbering has the property that node A is the ancestor of node B if and only if $A.Start < B.Start$ and $B.End < A.End$. In order to know the region numbers of the nodes associated with a tag, we attach the tag to the node that corresponds to the final location of the trace record of the tag. Thus, the region number that corresponds to the final node in the trace record of the tag is assigned to the tag. In order to get the time and location information for tag 6 in Fig. 7, we see the region number that corresponds to the final node in the trace record of tag 6. For the final node $C[7,8]$ for tag 6, its region number is $[13, 14]$. Therefore, we retrieve nodes that satisfy $Start \leq 13$ and $End \geq 14$. Such nodes are $A[2,3]$, $E[4,6]$, and $C[7,8]$. Therefore, we can retrieve the time information for tag 6 efficiently. There are more sophisticated region numbering schemes, however, the incorporation of them is straightforward.

6 STORING TRACE RECORDS IN AN RDBMS

We devise a relational schema to store RFID data based on the path encoding scheme and the region numbering scheme. The schema is shown in Fig. 8. $PATH_TABLE$, TAG_TABLE , and $TIME_TABLE$ are related to the trace records and $INFO_TABLE$ is related to the product information (e.g., product name, manufacturer, price).

$PATH_TABLE$ stores the path information by using the path encoding scheme in Section 5. In $PATH_TABLE$, the attributes $ELEMENT_ENC_1$, $ELEMENT_ENC_2$, ..., and $ELEMENT_ENC_m$ correspond to Element List Encoding Number, while the attributes $ORDER_ENC_1$, $ORDER_ENC_2$, ..., and $ORDER_ENC_m$ correspond to Order Encoding Number. If an RDBMS supports the data type whose maximum value is greater than Element List Encoding Number and Order Encoding Number of any trace record, $PATH_TABLE$ consists of only $ELEMENT_ENC_1$ and $ORDER_ENC_1$ with $PATH_ID$.

$TIME_TABLE$ stores the time information for trace records by using the region numbering scheme. In $TIME_TABLE$, $START$ and END are Start and End in the region numbering scheme. LOC is the location. $START_TIME$ and END_TIME correspond to the start time and end time in the time tree. TAG_TABLE has two identifiers for path and time information. $PATH_ID$ is the identifier for the path information and $(START, END)$ is the identifier for the time information. In addition, $INFO_ID$ is the identifier for

PATH_TABLE

PATH_ID	ELEMENT_ENC_1	...	ELEMENT_ENC_m	ORDER_ENC_1	...	ORDER_ENC_m
---------	---------------	-----	---------------	-------------	-----	-------------

TAG_TABLE

TAG_ID	PATH_ID	START	END	INFO_ID
--------	---------	-------	-----	---------

TIME_TABLE

START	END	LOC	START_TIME	END_TIME
-------	-----	-----	------------	----------

INFO_TABLE

INFO_ID	PRODUCT_NAME	MANUFACTURER	PRICE
---------	--------------	--------------	-------

Fig. 8. Relational schema to store RFID data.

INFO_TABLE. INFO_TABLE stores the information of products. In this paper, we do not focus on INFO_TABLE.

If a path is long, its Element List Encoding Number and Order Encoding Number will be large and commercial RDBMSs may not support the storing of such large numbers. To solve this problem, we store the path in several attributes by dividing it.

The notations used in this paper are defined as follows:

- *Type*: the data type for Element List Encoding Number and Order Encoding Number which an RDBMS supports (e.g., BIGINT, DECIMAL),
- Max_{Type} : the maximum value that the data type *Type* can have, and
- m : the maximum number of path segments in a path.

Consider the long path $L_1 \rightarrow L_2 \rightarrow \dots \rightarrow L_n$. We can divide the long path into several path segments s_1, \dots, s_t such that

-

$$\begin{aligned} s_1 &= L_1 \rightarrow L_2 \rightarrow \dots \rightarrow L_{i_1}, \\ s_2 &= L_{i_1+1} \rightarrow L_{i_1+2} \rightarrow \dots \rightarrow L_{i_2}, \\ &\dots, \\ s_t &= L_{i_{t-1}+1} \rightarrow L_{i_{t-1}+2} \rightarrow \dots \rightarrow L_{i_t} \\ &(1 \leq i_1 < i_2 < \dots < i_t = n) \end{aligned}$$

- For $u \in 1, 2, \dots, t$,

$$\begin{aligned} \prod_{j=i_{u-1}+1}^{i_u} Prime(L_j) &\leq Max_{Type}, \text{ and} \\ (\prod_{j=i_{u-1}+1}^{i_u} Prime(L_j)) \times Prime(L_{i_u+1}) &> Max_{Type}, \end{aligned}$$

with $i_0 = 0$ and $Prime(L_{n+1}) = \infty$.

We can easily find the path segments for the long path $L_1 \rightarrow L_2 \rightarrow \dots \rightarrow L_n$ which satisfy the above conditions. We do not consider Order Encoding Number when we divide the path since Order Encoding Number X for any path is less than Element Encoding Number E for the path by Theorem 2 (i.e., $0 \leq X < E$). However, elements in each path segment keep the order information in the original path. Example 3 shows how the path is split up.

Example 3. Assume that the path $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E \rightarrow F$, $p_1 = Prime(A) = 2$, $p_2 = Prime(B) = 3$, $p_3 = Prime(C) = 5$, $p_4 = Prime(D) = 7$, $p_5 = Prime(E) = 11$, $p_6 = Prime(F) = 13$, and $Max_{Type} = 250$. Since $\sum_{j=1}^6 p_j = 30030 > Max_{Type} = 250$, Element List Encoding Number for the path cannot be stored in one attribute with *Type*. We split up the path into path segments $s_1 : A \rightarrow B \rightarrow C \rightarrow D$, $s_2 : E \rightarrow F$ which satisfy the following inequalities. Note that the order of E in s_2 is not 1 but 5, and that of F is not 2 but 6.

$$\begin{aligned} \sum_{j=1}^4 p_j &= 210 < Max_{Type} \text{ and } (\sum_{j=1}^4 p_j) \times p_5 \\ &= 2310 > Max_{Type} \\ \sum_{j=5}^6 p_j &= 143 < Max_{Type} \text{ and } (\sum_{j=5}^6 p_j) \times p_7 \\ &= \infty > Max_{Type} \end{aligned}$$

```

Function store_trace_record(trace records tr)
begin
1: // m: the maximum number of path segments in a path
2: m := 0

3: for i=0; i<the number of trace records; i++
4: {
5:   <path_id, store_flag>:=constructTree(tree, tr[i])
6:   if store_flag == FALSE
7:   {
8:     <segment_list, segment_num> := divide_path(tr[i])
9:     if segment_num > m
10:    m := segment_num
11:    store segment_list for trace record tr[i] in PATH_FILE
        using the path encoding scheme
12:  }
13: store (tag identifier from tr[i], path_id) in TEMP_PATH_FILE
14: }

15: for i=0; i<the number of trace records; i++
16:   constructTimeTree(time_tree, tr[i])
17:   assign region numbers to nodes in time_tree

18: while traversing nodes in time_tree by the breath-first search
19: {
20:   store nodes of time_tree in TIME_FILE
21:   store (tag identifier, region numbers for node) in
        TEMP_TIME_FILE for all tags attached to node
22: }

23: create the schema according to m
24: perform bulk loading with PATH_FILE, TEMP_PATH_FILE,
        TIME_FILE, and TEMP_TIME_FILE
25: after joining TEMP_PATH_TABLE and TEMP_TIME_TABLE
        on TAG_ID, fill TAG_TABLE

end

```

Fig. 9. Algorithm to store trace records using the relational schema.

If the maximum number of path segments in a path is m , then PATH_TABLE is created, as shown in Fig. 8, which has m ELEMENT_ENC attributes and m ORDER_ENC attributes. If some path is divided into t path segments and t is less than m , we fill ELEMENT_ENC $_{(t+1)}$, ELEMENT_ENC $_{(t+2)}$, ..., and ELEMENT_ENC $_m$ with 1 and ORDER_ENC $_{(t+1)}$, ORDER_ENC $_{(t+2)}$, ... and ORDER_ENC $_m$ with 0.

Fig. 9 shows the algorithm to store trace records using the relational schema. As input for the algorithm, trace records are given. Trace records are translated into relational data. In the algorithm of Fig. 9, we pile up the relational data in a text file format and then store it into an RDBMS together with bulk loading in order to load data efficiently. PATH_FILE is the file for PATH_TABLE, TEMP_PATH_FILE for TEMP_PATH_TABLE, TIME_FILE for TIME_TABLE, and TEMP_TIME_FILE for TEMP_TIME_TABLE. In Line 2, we initialize m , the maximum number of path segments in a path. To fill PATH_TABLE, we construct a tree from the paths of the trace records tr by using $constructTree(Tree\ tree, TraceRecord\ tr)$ (Line 5). In $constructTree(Tree\ tree, TraceRecord\ tr)$, if there is no path for tr in $tree$, it inserts a new path into $tree$, returns $path_id$ of the new path and sets $store_flag$ to FALSE. Otherwise, it returns $path_id$ of the path and sets $store_flag$ to TRUE.

If $store_flag$ is FALSE, we insert Element List Encoding Number and Order Encoding Number for the path into PATH_FILE (Line 6-12). However, if the numbers are large and cannot be stored in one attribute, we divide the path using $divide_path(TraceRecord\ tr)$. The $divide_path(TraceRecord\ tr)$ divides one path into several path segments according to

PATH_TABLE			TAG_TABLE					TIME_TABLE				
PATH_ID	ELEMENT_ENC	ORDER_ENC	TAG_ID	PATH_ID	START	END	INFO_ID	START	END	LOC	START_TIME	END_TIME
1	30	23	tag1	1	3	4		1	20	A	2	3
2	42	17	tag2	1	3	4		2	7	B	5	7
3	110	13	tag3	1	3	4		8	11	B	7	8
4	22	13	tag4	2	5	6		12	15	E	4	6
5	14	9	tag5	2	9	10		16	17	D	4	5
			tag6	3	13	14		18	19	D	5	6
			tag7	3	13	14		3	4	C	8	9
			tag8	4	12	15		5	6	D	13	16
			tag9	5	16	17		9	10	D	14	18
			tag10	5	18	19		13	14	C	7	8

Fig. 10. Status of tables after storing trace records when m is 1.

Max_{Type} and returns a path segment list (i.e., *segment_list*) and the number of path segments in the path (i.e., *segment_num*). Each path segment is stored in *PATH_FILE* by using the path encoding scheme with *path_id*. Also, for the creation of the schema, we keep the maximum number of path segments, m (Line 9-10). In Line 13, we insert *tag identifier* and *path_id* into *TEMP_PATH_FILE*. Note that *TEMP_PATH_TABLE* is used to fill *TAG_TABLE* later.

To fill *TIME_TABLE*, we construct the time tree from the trace records *tr* by using *constructTimeTree(TimeTree time_tree, TraceRecord tr)* (Line 16). In *constructTimeTree(TimeTree time_tree, TraceRecord tr)*, if there is no *tr* in *time_tree*, it inserts the new trace record into *time_tree*. After the construction of *time_tree*, we assign region numbers to the nodes in *time_tree* (Line 17). Then, we store (START, END, LOC, START_TIME, END_TIME) in *TIME_FILE* (Line 20) by traversing nodes in *time_tree* by using the breath-first search. If there are tags attached to a node, we store (tag identifier, region number for the node) in *TEMP_TIME_FILE* (Line 21). We create the schema according to the maximum number of path segments m . We then perform bulk loading with *PATH_FILE*, *TEMP_PATH_FILE*, *TIME_FILE*, and *TEMP_TIME_FILE*. Finally, to fill *TAG_TABLE*, we join *TEMP_PATH_TABLE* and *TEMP_TIME_TABLE* on *TAG_ID* (Line 25).

Fig. 10 shows the status of tables after storing trace records in Fig. 5 by the algorithm in Fig. 9. In Fig. 10, we assume that all Element List Encoding Numbers can be stored in one attribute since Max_{Type} is large enough. Since there are five different paths for trace records in Fig. 5, Element List Encoding Numbers and Order Encoding Numbers of five paths are stored in *PATH_TABLE*. Note that the path identifiers for $A \rightarrow B \rightarrow C$, $A \rightarrow B \rightarrow D$, $A \rightarrow E \rightarrow C$, $A \rightarrow E$, and $A \rightarrow D$ are 1, 2, 3, 4, and 5, respectively. As shown in Fig. 7, the time tree for trace records has 10 nodes. The information for 10 nodes is stored in *TIME_TABLE*. Also, the path identifier and the time identifier for 10 tags are stored in *TAG_TABLE*.

If we cannot store Element List Encoding Number in one attribute, we should split up the paths as mentioned above. Fig. 11 shows the status of *PATH_TABLE* after storing trace records in Fig. 5 when Max_{Type} is 30. In this case, m is 2. Also, since *TAG_TABLE* and *TIME_TABLE* are not affected by m , they are the same as those in Fig. 10. Thus, we do not show them in Fig. 11. In Fig. 11, the path $A \rightarrow B \rightarrow D$ (*path_id* 2) and the path $A \rightarrow E \rightarrow C$ (*path_id* 3) are divided into two path segments. Therefore, those

paths are stored with two *ELEMENT_ENC* attributes and two *ORDER_ENC* attributes. For example, $A \rightarrow B \rightarrow D$ has two path segments, $A \rightarrow B$ and D . $A \rightarrow B$ is encoded as *ELEMENT_ENC_1* = 2×3 and *ORDER_ENC_1* = 5 and D as *ELEMENT_ENC_2* = 7 and *ORDER_ENC_2* = 3. Other paths are stored with one *ELEMENT_ENC* attribute and one *ORDER_ENC* attribute and their remaining attributes are filled with 1 or 0.

7 QUERY TRANSLATION

Based on the relational schema in Section 6, we provide a method to process tracking queries and path-oriented queries. Since we store RFID data by using an RDBMS, we translate tracking queries and path-oriented queries into SQL queries. Then, we can easily process the queries by executing the SQL queries. We will first deal with the translation when m is 1 in Section 7.1 and then deal with the translation when m is greater than 1 in Section 7.2.

7.1 Basic Query Translation

In this section, we assume that all Element List Encoding Numbers can be stored in one attribute. That is, m is 1 and *PATH_TABLE* consists of three attributes: *path_id*, *ELEMENT_ENC_1*, and *ORDER_ENC_1*. Therefore, for convenience, we use *ELEMENT_ENC* and *ORDER_ENC* instead of *ELEMENT_ENC_1* and *ORDER_ENC_1*, respectively.

7.1.1 Tracking Query

We can process a tracking query efficiently by using the relational schema in Section 6. To trace a tag, we get Element List Encoding Number and Order Encoding Number that correspond to the tag. To get them, we join *PATH_TABLE* and *TAG_TABLE*. Fig. 12 shows the SQL query to get Element List Encoding Number and Order

PATH_TABLE				
PATH_ID	ELEMENT_ENC_1	ELEMENT_ENC_2	ORDER_ENC_1	ORDER_ENC_2
1	30	1	23	0
2	6	7	5	3
3	22	5	13	3
4	22	1	13	0
5	14	1	9	0

Fig. 11. Status of tables after storing trace records when m is 2.

```
SELECT P.ELEMENT_ENC, P.ORDER_ENC
FROM PATH_TABLE P, TAG_TABLE T
WHERE T.TAG_ID = my_tag_id AND T.PATH_ID = P.PATH_ID
```

Fig. 12. SQL query for the tracking query.

```
SELECT T.TAG_ID
FROM PATH_TABLE P, TAG_TABLE T
WHERE MOD(P.ELEMENT_ENC, pA*pB*pC) = 0 AND
MOD(P.ORDER_ENC, pA) < MOD(P.ORDER_ENC, pB) AND
MOD(P.ORDER_ENC, pB) + 1 = MOD(P.ORDER_ENC, pC)
AND P.PATH_ID = T.PATH_ID
```

Fig. 13. SQL query for <>//A//B/C>.

```
SELECT T.TAG_ID
FROM PATH_TABLE P, TAG_TABLE T, TIME_TABLE S
WHERE MOD(P.ELEMENT_ENC, pA*pB*pC) = 0 AND
MOD(P.ORDER_ENC, pA) < MOD(P.ORDER_ENC, pB) AND
MOD(P.ORDER_ENC, pB) + 1 = MOD(P.ORDER_ENC, pC)
AND P.PATH_ID = T.PATH_ID AND S.LOC = 'B' AND
S.START <= T.START AND T.END <= S.END AND
(S.END_TIME - S.START_TIME) < 10
```

Fig. 14. SQL query for <>//A//B[(EndTime-StartTime)<10]/C>.

Encoding Number that correspond to the tag with TagID = my_tag_id .

In order to know the locations in the flow of the tag, we factorize Element List Encoding Number. We then order the prime number factors P_1, \dots, P_n by computing *Order Encoding Number mod P_i* . We finally transform the prime number into the location name that corresponds to it.

Though the SQL query in Fig. 12 has the join between PATH_TABLE and TAG_TABLE, it does not have too much time to execute the query since the query has the selection predicate for TAG_TABLE (T.TAG_ID = my_tag_id) and there is only one tuple that satisfies the predicate. Therefore, we can process tracking queries efficiently.

7.1.2 Path-Oriented Retrieval Query

Although the path condition in a path-oriented retrieval query has many ancestor-descendant relationships, we can easily find paths that satisfy the path condition by checking mathematical conditions. Therefore, we can process path-oriented retrieval queries efficiently with our relational schema.

By Theorem 1, the path contains locations L_1, L_2, \dots, L_k if and only if $ELEMENT_ENC \bmod (Prime(L_1) \times Prime(L_2) \times \dots \times Prime(L_k)) = 0$. Therefore, if the path condition contains locations L_1, L_2, \dots, L_k , we insert the element membership condition $ELEMENT_ENC \bmod (Prime(L_1) \times Prime(L_2) \times \dots \times Prime(L_k)) = 0$ into the where clause in the SQL query. To determine the ancestor-descendant relationship or the parent-child relationship, we use Order Encoding Number. Consider two adjacent steps in the path condition and let locations of the two steps be L_a and L_b . If L_a and L_b are in the ancestor-descendant relationship (i.e., $L_a // L_b$), we insert $ELEMENT_ENC \bmod Prime(L_a) < ELEMENT_ENC \bmod Prime(L_b)$ into the where clause in the SQL query. If L_a and L_b are in the parent-child relationship (i.e., L_a / L_b), we insert $ELEMENT_ENC \bmod$

```
SELECT T.TAG_ID
FROM PATH_TABLE P, TAG_TABLE T, INFO_TABLE I
WHERE MOD(P.ELEMENT_ENC, pA*pB*pC) = 0 AND
MOD(P.ORDER_ENC, pA) < MOD(P.ORDER_ENC, pB) AND
MOD(P.ORDER_ENC, pB) + 1 = MOD(P.ORDER_ENC, pC)
AND P.PATH_ID = T.PATH_ID AND
I.PRODUCT_NAME = 'laptop' AND I.INFO_ID = T.INFO_ID
```

Fig. 15. SQL query for <>//A//B/C, PRODUCT_NAME = 'laptop'>.

```
SELECT COUNT(*)
FROM PATH_TABLE P, TAG_TABLE T
WHERE MOD(P.ELEMENT_ENC, pA*pB*pC) = 0 AND
MOD(P.ORDER_ENC, pA) < MOD(P.ORDER_ENC, pB) AND
MOD(P.ORDER_ENC, pB) + 1 = MOD(P.ORDER_ENC, pC)
AND P.PATH_ID = T.PATH_ID
```

Fig. 16. SQL query for <COUNT(), //A//B/C>.

```
SELECT AVG(S.START_TIME)
FROM PATH_TABLE P, TAG_TABLE T, TIME_TABLE S
WHERE MOD(P.ELEMENT_ENC, pA*pB*pC) = 0 AND
MOD(P.ORDER_ENC, pA) < MOD(P.ORDER_ENC, pB) AND
MOD(P.ORDER_ENC, pB) + 1 = MOD(P.ORDER_ENC, pC)
AND P.PATH_ID = T.PATH_ID AND S.LOC = 'B' AND
S.START <= T.START AND T.END <= S.END
```

Fig. 17. SQL query for <AVG(B.StartTime), //A//B/C>.

```
SELECT AVG(S2.END_TIME - S1.START_TIME)
FROM PATH_TABLE P, TAG_TABLE T, TIME_TABLE S1, TIME_TABLE S2
WHERE MOD(P.ELEMENT_ENC, pA*pB*pC) = 0 AND
MOD(P.ORDER_ENC, pA) < MOD(P.ORDER_ENC, pB) AND
MOD(P.ORDER_ENC, pB) + 1 = MOD(P.ORDER_ENC, pC) AND
P.PATH_ID = T.PATH_ID AND S1.LOC = 'A' AND
S1.START <= T.START AND T.END <= S1.END AND S2.LOC = 'C' AND
S2.START <= T.START AND T.END <= S2.END
```

Fig. 18. SQL query for <AVG(C.EndTime - A.StartTime), //A//B/C>.

$Prime(L_a) + 1 = ELEMENT_ENC \bmod Prime(L_b)$ into the where clause.

After finding the paths that satisfy the path condition, we join PATH_TABLE and TAG_TABLE on TAG_ID to get tags. Fig. 13 shows the SQL query that corresponds to the path-oriented retrieval query <>//A//B/C>. In Figs. 13, 14, 15, 16, 17, and 18, pA , pB , and pC denote $Prime(A)$, $Prime(B)$, and $Prime(C)$, respectively.

Path-oriented retrieval queries may have time conditions. If the queries have time conditions, we join TAG_TABLE and TIME_TABLE. We can retrieve the time information efficiently by using the property of the region numbering scheme. We insert the following statement into the where clause in the SQL query for time conditions:

```
TIME_TABLE.LOC = 'location name' AND
TIME_TABLE.START <= TAG_TABLE.START AND
TAG_TABLE.END <= TIME_TABLE.END AND
Time Conditions in the Step
```

Fig. 14 shows the SQL query that corresponds to the path-oriented retrieval query <>//A//B[(EndTime-StartTime)<10]/C>.

Path-oriented retrieval queries may also have product information conditions, such as PRODUCT_NAME = 'laptop'. To process such queries, we first perform the selection of

INFO_TABLE for product information conditions. We then join INFO_TABLE and TAG_TABLE on INFO_ID. Fig. 15 shows the SQL query that corresponds to the path-oriented retrieval query $\langle //A//B/C, PRODUCT_NAME = 'laptop' \rangle$.

7.1.3 Path-Oriented Aggregate Query

Since path-oriented aggregate queries have aggregate functions, we add an aggregate function in the select clause of the SQL query. Fig. 16 shows the SQL query that corresponds to the path-oriented aggregate query $\langle COUNT(), //A//B/C \rangle$.

In the case of aggregate functions that need time attributes as arguments, we join TAG_TABLE and TIME_TABLE to get the time attributes since PATH_TABLE does not have the time information. Fig. 17 shows the SQL query that corresponds to the path-oriented aggregate query $\langle AVG(B.StartTime), //A//B/C \rangle$.

Consider the query $\langle AVG(C.EndTime - A.StartTime), //A//B/C \rangle$. The query has the aggregate function that has two time attributes as the argument. In this case, we join one TAG_TABLE and two TIME_TABLES. Fig. 18 shows the SQL query that corresponds to the path-oriented aggregate query $\langle AVG(C.EndTime - A.StartTime), //A//B/C \rangle$.

7.2 Advanced Query Translation

In this section, we assume that Element List Encoding Number cannot be stored in one attribute. Although we can store Element List Encoding Number for a long path in an RDBMS by dividing the path, it is difficult to find the paths that satisfy the path condition using the SQL query due to the computation overflow. For example, to find some path from PATH_TABLE in Fig. 11, we should multiply ELEMENT_ENC_1 by ELEMENT_ENC_2. However, since Max_{Type} is 30, the overflow will occur in an RDBMS during the multiplication. The SQL query may not run, or the result may be wrong. Therefore, we need to extend the translation for path conditions when m is greater than 1.

When a path condition contains locations L_1, L_2, \dots, L_k , we should insert the element membership condition

$$(ELEMENT_ENC_1 \times ELEMENT_ENC_2 \times \dots \times ELEMENT_ENC_m) \bmod (Prime(L_1) \times Prime(L_2) \times \dots \times Prime(L_k)) = 0$$

into the where clause in the SQL query in order to check whether the path contains locations L_1, L_2, \dots, L_k . To avoid the computation overflow during the execution of the condition, we change the condition into another equivalent condition by using Lemmas 1 and 2.

Lemma 1. Let E_1, E_2 , and P be any natural numbers. Then,

$$(E_1 E_2) \bmod P = (E_1 \bmod P)(E_2 \bmod P) \bmod P.$$

Proof. This is easily derived from the basic number theory. \square

Lemma 2. Let E, P_1 , and P_2 be any natural numbers and P_1 and P_2 be relatively prime numbers. Then,

$$E \bmod (P_1 P_2) = 0 \text{ if and only if } E \bmod P_1 = 0 \text{ and } E \bmod P_2 = 0.$$

Proof. This is easily derived from the basic number theory. \square

```

Function element_translation(element list  $L_1, L_2, \dots, L_k$  in the path condition)
begin
1: product := 1
2: condition := null
3: for i=1; i<=k; i++
4: {
5:   if (product * Prime( $L_i$ ))2 > MaxType
6:   {
7:     if condition == null
8:       // m: the maximum number of path segments in a path
9:       condition := sub_element_translation(product, m) + " = 0 "
10:    else
11:    {
12:      // m: the maximum number of path segments in a path
13:      condition := condition + " and " + sub_element_translation(product, m) + " = 0 "
14:    }
15:    product := 1
16:  }
17:  product := product * Prime( $L_i$ )
18: }
19: if condition == null
20:   condition := sub_element_translation(product, m) + " = 0 "
21: else
22: {
23:   condition := condition + " and " + sub_element_translation(product, m) + " = 0 "
24: }
25: return condition
end

Function sub_element_translation(product p, index i)
begin
1: if i==1
2:   return "ELEMENT_ENC_" + i + " mod " + p
3: else
4:   return "(" + sub_element_translation(p, i-1) + ") (ELEMENT_ENC_" + i +
      " mod " + p + ") mod " + p
end

```

Fig. 19. Algorithm to Translate the Element Membership Condition into an SQL Query.

Suppose that the path condition in the query contains L_1, L_2, \dots, L_k . Let p_i be $Prime(L_i)$, P be $p_1 p_2 \dots p_k$, and E_i be $ELEMENT_ENC_i$. To know whether the path contains locations L_1, L_2, \dots, L_k , we should check the condition $(E_1 E_2 \dots E_m) \bmod P = 0$. To avoid the overflow during the computation of the condition, we change it into another equivalent condition using Lemma 1, which can be computed without an overflow. Consider the simple case that m is 2. By Lemma 1, $(E_1 E_2) \bmod P$ is equal to $(E_1 \bmod P)(E_2 \bmod P) \bmod P$. Therefore, we check $(E_1 \bmod P)(E_2 \bmod P) \bmod P = 0$ instead of $(E_1 E_2) \bmod P = 0$. If $P^2 \leq Max_{Type}$, the computation of $(E_1 \bmod P)(E_2 \bmod P) \bmod P = 0$ is performed without an overflow since $A \bmod P < P$, for any natural number A .

If m is greater than 2, we apply Lemma 1 recursively. The detailed algorithm is shown in Fig. 19. However, the translation using Lemma 1 is limited since $P^2 \leq Max_{Type}$ should be satisfied. If $P^2 \leq Max_{Type}$ is not satisfied, we can apply Lemma 2 before using Lemma 1. To apply Lemma 2, we first factorize P into P_1, P_2, \dots, P_t such that $P = P_1 P_2 \dots P_t$ and $P_i^2 \leq Max_{Type}$ for $i = 1, 2, \dots, t$. If P is factorized into P_1 and P_2 , we can change the condition $(E_1 E_2) \bmod P = 0$ into two conditions $(E_1 E_2) \bmod P_1 = 0$ and $(E_1 E_2) \bmod P_2 = 0$ by using Lemma 2. Also, by using Lemma 1, we can change $(E_1 E_2) \bmod P_1 = 0$ and $(E_1 E_2) \bmod P_2 = 0$ into $(E_1 \bmod P_1)(E_2 \bmod P_1) \bmod P_1 = 0$ and $(E_1 \bmod P_2)(E_2 \bmod P_2) \bmod P_2 = 0$. An RDBMS can evaluate $(E_1 \bmod P_1)(E_2 \bmod P_1) \bmod P_1 = 0$ and $(E_1 \bmod P_2)(E_2 \bmod P_2) \bmod P_2 = 0$ without a computation overflow.

The algorithm in Fig. 19 summarizes the above process for checking whether the path contains locations

```

SELECT T.TAG_ID
FROM PATH_TABLE P, TAG_TABLE T
WHERE
MOD(MOD(P.ELEMENT_ENC_1, pA*pB*pC) *
MOD(P.ELEMENT_ENC_2, pA*pB*pC) = 0 AND
MOD(P.ORDER_ENC_1, pA) * FLOOR(1 - (P.ELEMENT_ENC_1 mod pA) / pA) +
MOD(P.ORDER_ENC_2, pA) * FLOOR(1 - (P.ELEMENT_ENC_2 mod pA) / pA)
< MOD(P.ORDER_ENC_1, pB) * FLOOR(1 - (P.ELEMENT_ENC_1 mod pB) / pB) +
MOD(P.ORDER_ENC_2, pB) * FLOOR(1 - (P.ELEMENT_ENC_2 mod pB) / pB) AND
MOD(P.ORDER_ENC_1, pB) * FLOOR(1 - (P.ELEMENT_ENC_1 mod pB) / pB) +
MOD(P.ORDER_ENC_2, pB) * FLOOR(1 - (P.ELEMENT_ENC_2 mod pB) / pB) + 1 =
MOD(P.ORDER_ENC_1, pC) * FLOOR(1 - (P.ELEMENT_ENC_1 mod pC) / pC) +
MOD(P.ORDER_ENC_2, pC) * FLOOR(1 - (P.ELEMENT_ENC_2 mod pC) / pC)

```

Fig. 20. SQL query for $\langle //A//B/C \rangle$ when m is 2.

L_1, L_2, \dots, L_k . As input for the algorithm *element_translation()*, the element list L_1, L_2, \dots, L_k in the path condition is given. In Line 1-2, *product* and *condition* are initialized. If the square of the product of prime numbers for the elements in the path is greater than Max_{Type} , we call *sub_element_translation(product, m)* and reinitialize *product*. *sub_element_translation(product, m)* translates $(E_1 E_2 E_3 \dots E_m) \bmod P$ into the equivalent condition, which does not cause the computation overflow.

To determine the ancestor-descendant relationship or the parent-child relationship, we use Order Encoding Number. When m is 1, we compute $ORDER_ENC_1 \bmod p_j$ for the order of location L_j . However, when m is greater than 1, the order information of each location in a path is separately stored in attributes $ORDER_ENC_1, ORDER_ENC_2, \dots$, and $ORDER_ENC_m$. We use Theorem 3 to compute the order when m is greater than 1.

Theorem 3. Let O be Order Encoding Number in an original path before dividing the path and O_i be $ORDER_ENC_i$ that corresponds to the i th path segment after dividing the path. For each location L_j in the path, the following equation is satisfied:

$$O \bmod p_j = \sum_{i=1}^m (O_i \bmod p_j) \lfloor 1 - \frac{E_i \bmod p_j}{p_j} \rfloor.$$

Proof. If L_j is contained in the i th path segment, $(O_i \bmod p_j) \lfloor 1 - \frac{E_i \bmod p_j}{p_j} \rfloor = (O_i \bmod p_j) = (O \bmod p_j)$ since $E_i \bmod p_j = 0$. Otherwise, $(O_i \bmod p_j) \lfloor 1 - \frac{E_i \bmod p_j}{p_j} \rfloor = 0$ since $E_i \bmod p_j \neq 0$ and $E_i \bmod p_j < p_j$. Therefore, $\sum_{j=1}^m (O_i \bmod p_j) \lfloor 1 - \frac{E_i \bmod p_j}{p_j} \rfloor$ computes the order information of location L_j in m $ORDER_ENC$ attributes. That is, $O \bmod p_j = \sum_{i=1}^m (O_i \bmod p_j) \lfloor 1 - \frac{E_i \bmod p_j}{p_j} \rfloor$. \square

If we use $\sum_{i=1}^m (O_i \bmod p_j) \lfloor 1 - \frac{E_i \bmod p_j}{p_j} \rfloor$ in Theorem 3, we can translate the order condition into an SQL query. We use the *FLOOR* operator in SQL queries for $\lfloor \cdot \rfloor$.

Fig. 20 shows the SQL query that corresponds to the path-oriented retrieval query $\langle //A//B/C \rangle$ when m is 2. We omit other translation examples since the translation, except for the path information, is similar to that of Section 7.1.

8 EXPERIMENTS

In order to validate our approach, we conduct experimental evaluations for various queries.

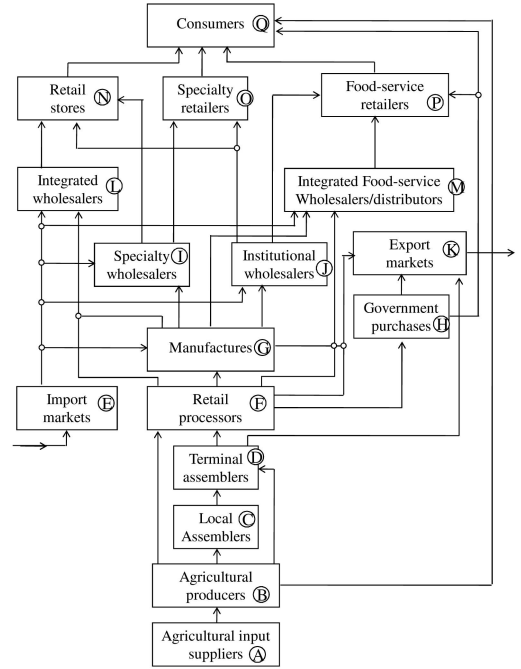


Fig. 21. Food distribution in the United States [9].

8.1 Experimental Environment

We experiment on a Pentium 2.53 GHz with 1 GB main memory using Java. Since there is no well-known RFID data set, we generate synthetic data based on a real-life example and formulate 15 queries (1 tracking query, 4 path-oriented retrieval queries, 10 path-oriented aggregate queries). The query performance is measured by processing the queries three times and averaging the execution time. As a comparison system, RFID-Cuboid [14] is used. For fairness, we implement RFID-Cuboid on an RDBMS to support tracking queries and path-oriented queries.

8.1.1 Data Set

Since we use stay records instead of raw RFID data, we generate stay records. To reflect real environments, we generate data sets based on a real-life example in Fig. 21. Fig. 21 shows the food distribution [9] in the United States. In the food distribution, the minimum length of paths is 3 and the maximum length is 9. The products are distributed from import markets and agricultural input suppliers. Since it is difficult to find detailed distribution information of each product from a real-life example, we assume that the products are distributed equally when they move to the next location if there are many next locations.

Initially, products are generated at starting nodes such as import markets and agricultural input suppliers. Then, they move to the next location in groups according to the links of Fig. 21. As products move together in groups in many RFID applications, we generate stay records with the grouping factor, which is the number of data generated at a starting node simultaneously. We consider two kinds of data in data generation: GData (grouping factor: 10,000) and IData (grouping factor: 1,000). Although we set the grouping factor of GData to 10,000 and that of IData to 1,000, the values (i.e., 10,000, 1,000) are not important. The key point in

Query Number	Query	Query Number	Query
Query 1	TagID = 1	Query 9	<COUNT(), //D0//F0[StartTime<2000]//I0//O0>
Query 2	<//D0>	Query 10	<AVG(D0.EndTime-D0.StartTime), //D0//F0//I0//O0>
Query 3	</A0/B0/C0/D0/F0>	Query 11	<AVG(D0.EndTime-D0.StartTime), //D0//F0//I0//O0>
Query 4	<//D0//F0//I0//O0>	Query 12	<AVG(D0.StartTime), //D0>
Query 5	<//D0//F0[StartTime<2000]//I0//O0>	Query 13	<MIN(D0.EndTime-D0.StartTime), //D0//F0//I0//O0>
Query 6	<COUNT(), //D0>	Query 14	<MIN(D0.EndTime-D0.StartTime), //D0//F0//I0//O0>
Query 7	<COUNT(), //D0//F0//I0//O0>	Query 15	<MIN(D0.StartTime), //D0>
Query 8	<COUNT(), //D0//F0//I0//O0>		

Fig. 22. Query set.

generating GData and IData is that products in GData move in larger groups compared to those in IData. By generating GData and IData, we analyze how the query performance is affected by the grouping factor. We set the number of instances for each location to 2 since we consider all possible paths and the number of paths increases exponentially.

Here, 6×10^6 , 12×10^6 , 18×10^6 , 24×10^6 , and 30×10^6 stay records are generated for both GData and IData. The DECIMAL type in an RDBMS is used to store Element List Encoding Number and Order Encoding Number.

8.1.2 Query Set

Fifteen queries are formulated to test various features. Query 1 is a tracking query, Queries 2-5 are path-oriented retrieval queries, and Queries 6-15 are path-oriented aggregate queries. The queries are shown in Fig. 22. Note that we represent the location name in queries as the character in the circle in Fig. 21 and the instance identifier. Query 1 tests the performance of the tracking query. Queries 2, 6, 12, and 15 evaluate the performance of the path condition with a single location, while other queries with multiple locations. In particular, Queries 5 and 9 have the time condition.

8.2 Experimental Results

In experimental results, if the query is not finished within 2 hours, we denote it by "X" in the graph.

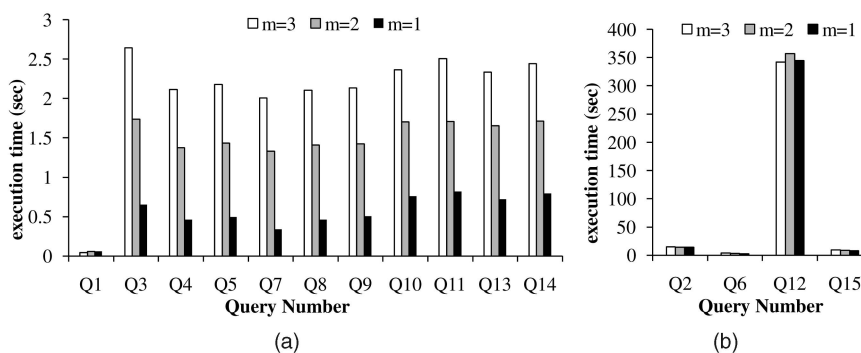
Fig. 23 shows the query performance in IData with 3×10^7 as m , the maximum number of path segments in a path, changes. Fig. 23a shows the query performance when the execution time is less than 3 seconds and Fig. 23b shows the performance when the time is more than 3 seconds. As m

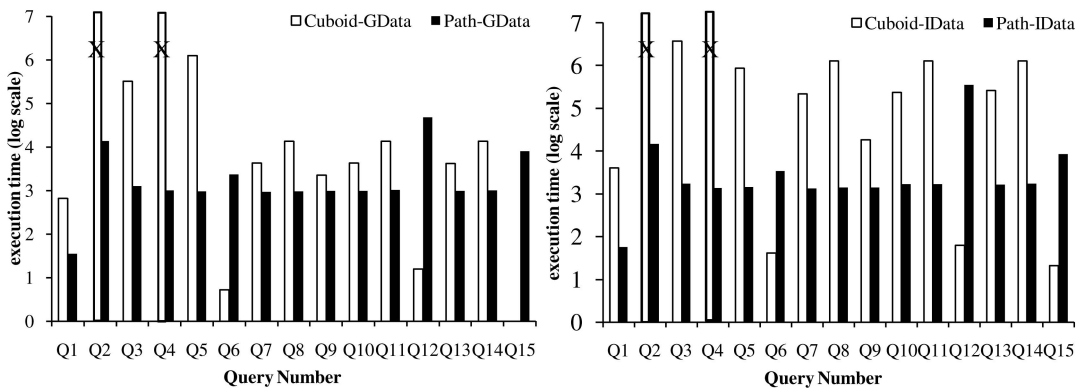
increases, the computation overhead becomes generally large since we need the additional computation for segmented paths. However, as can be seen in Fig. 23b, the performance sometimes is not affected much by m . This is partly because the RDBMS does not properly optimize the translated SQL query as the query is quite complex. While a lower value of m is favorable for query performance, a higher value must be used to avoid the overflow of the DECIMAL type. This is especially necessary for RDBMSs that do not provide a large DECIMAL type. Therefore, we set m to 2 for the following experiments.

Fig. 24 shows the query performance for 15 queries. Figs. 24a and 24b show the performance with 3×10^7 tuples in GData and IData, respectively. Path denotes our approach and Cuboid denotes RFID-Cuboid in the figures of this section. Since the performance gap between our approach and RFID-Cuboid is wide, we use the logarithmic scale (base 10) for the execution time in Fig. 24.

In GData and IData, our approach is better than RFID-Cuboid for most queries in terms of the execution time (except Queries 6, 12, and 15). We can also observe that the performance gap between our approach and RFID-Cuboid in IData is larger than that in GData. Since the shapes of the graphs for the execution time for different sizes are similar, the execution time for only GData with 3×10^7 and IData with 3×10^7 is shown in Fig. 24.

Fig. 25, 26, and 27 show the query performance according to the number of stay records. Fig. 25a shows the performance for a tracking query. For the tracking query, our approach is faster than RFID-Cuboid in both GData and IData. We can also observe that our approach is much faster than RFID-Cuboid in IData. While our

Fig. 23. Experiment with respect to m . (a) When the execution time is less than 3 seconds. (b) When the execution time is more than 3 seconds.



* In Cuboid-GData, the execution time for Query 15 is 0 in terms of microseconds. Therefore, we cannot represent it with the logarithmic scale.

(a)

(b)

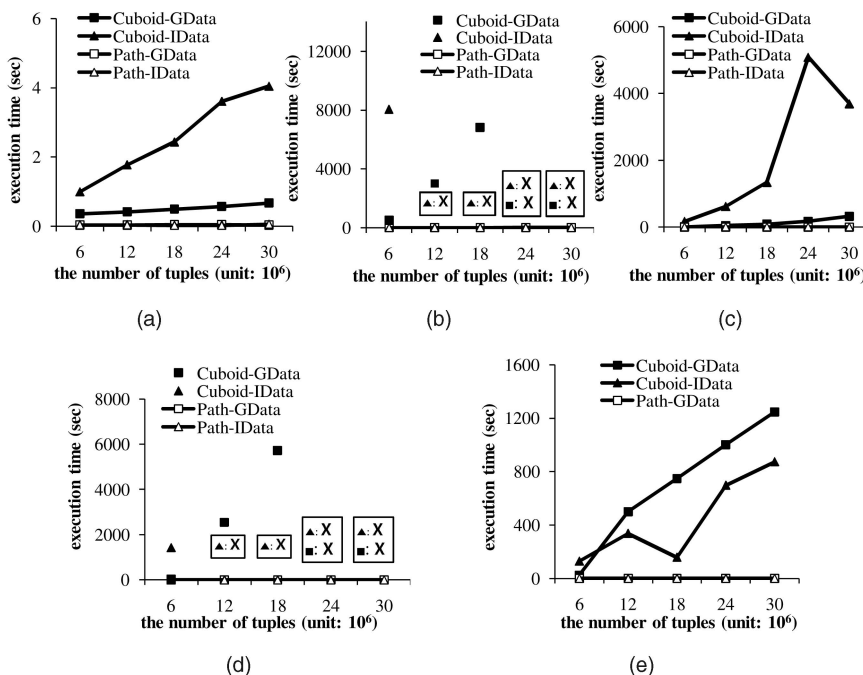
Fig. 24. Execution time for 15 queries. (a) GData. (b) IData.

approach finds only Element List Encoding Number and Order Encoding Number for the given tag identifier, RFID-Cuboid scans STAY_TABLE (the table for stay records). Therefore, in IData, RFID-Cuboid has much worse performance than our approach compared to GData since the number of tuples of STAY_TABLE in IData is more than that in GData.

The performance of Queries 2, 3, and 4 is shown in Figs. 25b, 25c, and 25d. However, most cases of RFID-Cuboid do not finish within 2 hours. They are denoted by X in a rectangle. For example, $\blacksquare: X$ in a rectangle on the horizontal axis 24 means that the \blacksquare approach does not finish within 2 hours when the number of tuples is 24×10^6 . These queries are path-oriented retrieval queries. To process path-oriented retrieval queries, RFID-Cuboid

joins STAY_TABLE and MAP_TABLE. MAP_TABLE contains the mapping from GID to TAG_ID in RFID-Cuboid. Since RFID-Cuboid uses the prefix encoding scheme, it needs string comparisons to process the join between STAY_TABLE and MAP_TABLE. Therefore, for path-oriented retrieval queries, our approach has a much better performance than RFID-Cuboid.

The performance of Query 5 is shown in Fig. 25e. Although RFID-Cuboid in GData shows a much better performance than that in IData for the cases of Queries 2, 3, and 4, the performance gap between GData and IData in RFID-Cuboid does not have a big difference in Fig. 25e. On the contrary, the performance in IData is better than that in GData. Query 5 has the time condition (StartTime<2,000).



(a)

(b)

(c)

(d)

(e)

Fig. 25. Execution time for Queries 1, 2, 3, 4, and 5. (a) Query 1. (b) Query 2. (c) Query 3. (d) Query 4. (e) Query 5.

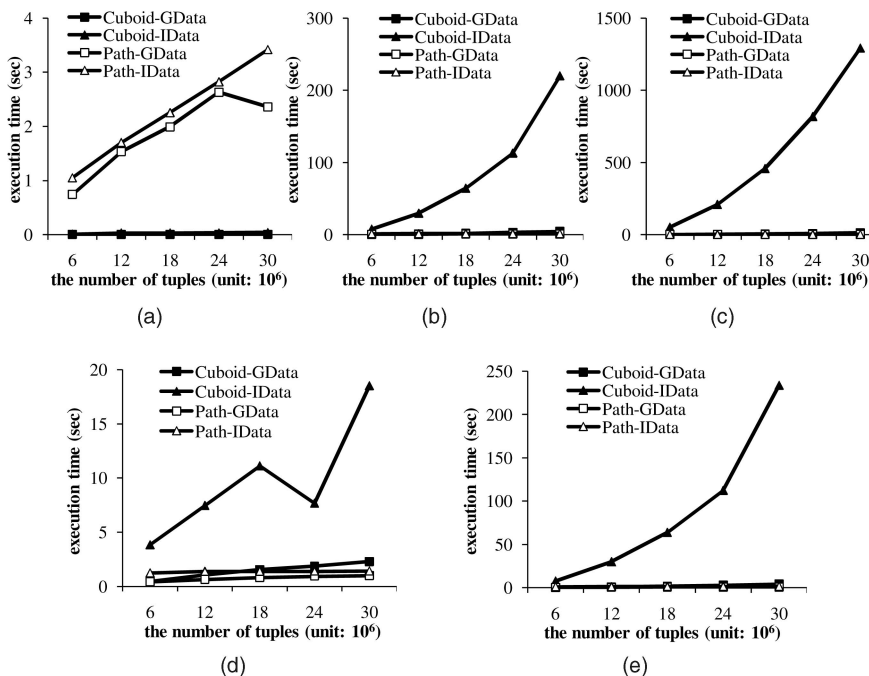


Fig. 26. Execution time for Queries 6, 7, 8, 9, and 10. (a) Query 6. (b) Query 7. (c) Query 8. (d) Query 9. (e) Query 10.

Therefore, in IData, many tuples of STAY_TABLE are removed by the condition and tuples to join are reduced significantly. Also, the number of results in IData is less than that in GData. Therefore, the performance gap between IData and GData for RFID-Cuboid is small for Query 5.

In Queries 6, 12, and 15 (Fig. 26a, Fig. 27b, and Fig. 27e), RFID-Cuboid is better than our approach. The path condition in Queries 6, 12, and 15 has only one location.

Since RFID-Cuboid focuses on groups in which products move together, RFID-Cuboid is efficient in the case of getting information at one location. However, in supply chain management, it is important to analyze the object transition. For queries related to the object transition (Queries 3, 4, 5, 7, 8, 9, 10, 11, 13, and 14), our approach is superior to RFID-Cuboid. However, for Query 2, our approach has a better performance than RFID-Cuboid,

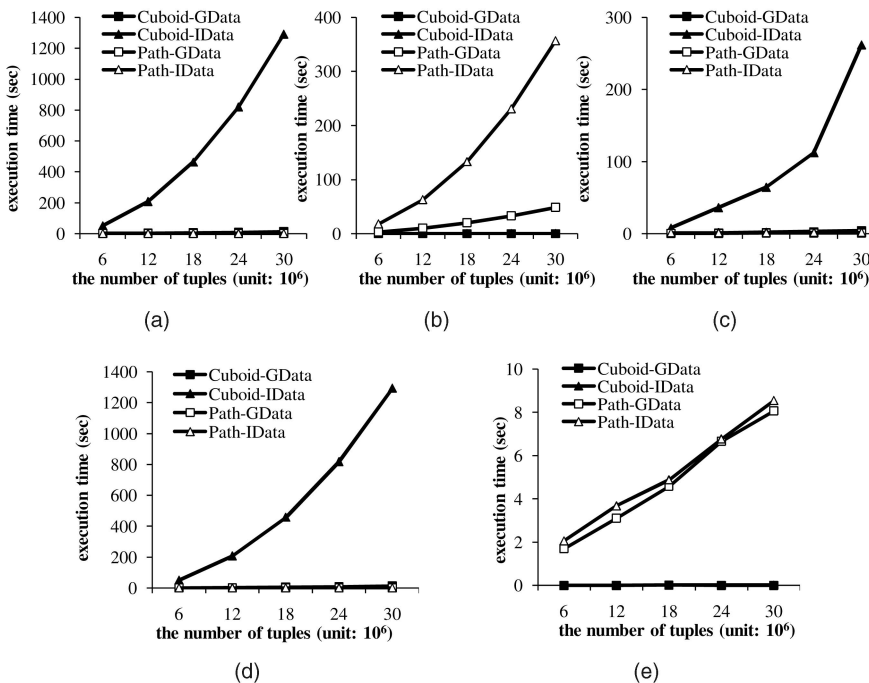


Fig. 27. Execution time for Queries 11, 12, 13, 14, and 15. (a) Query 11. (b) Query 12. (c) Query 13. (d) Query 14. (e) Query 15.

although Query 2 is used to get information at one location. This is because RFID-Cuboid uses the string comparison to get tags.

Consider the query performance of Queries 6, 12, and 15 (Fig. 26a, Fig. 27b, and Fig. 27e) versus that of Query 7, 8, 10, 11, 13, and 14 (Fig. 26b, Fig. 26c, Fig. 26e, Fig. 27a, Fig. 27c, and Fig. 27d). The path condition in Queries 6, 12, and 15 has only one location while the path condition in Queries 7, 8, 10, 11, 13, and 14 has multiple locations. Since our approach uses Element List Encoding Number and Order Encoding Number, our approach can easily find paths, although the path condition has many ancestor-descendant relationships. Therefore, in Queries 7, 8, 10, 11, 13, and 14, our approach is better than RFID-Cuboid.

Since Query 9 has the time information, our approach joins TAG_TABLE and TIME_TABLE. Therefore, our approach has a little better performance than RFID-Cuboid for GData, as shown in Fig. 26d.

Consequently, our approach is superior to RFID-Cuboid in most cases. Also, our approach is less sensitive than RFID-Cuboid for the grouping factor. Therefore, our approach can be used in a wide range of applications. For the path-oriented aggregate query whose path condition has only one location, our approach may be worse than RFID-Cuboid. However, since it is important to analyze the object transition in supply chain management, our approach will be more useful.

9 CONCLUSION

We expect that RFID technology will revolutionize supply chain management. In supply chain management, a large amount of RFID data is generated. However, since RFID data have flow information that is different from the traditional data, it is difficult to store data and process queries. Therefore, we propose an efficient storage scheme and efficient query processing for supply chain management.

The proposed approach in this paper is very efficient in tracking and analyzing the flow of products. We believe that it is useful for real-life logistic applications after conducting experiments with data sets based on a real-life logistic example. In addition, we devise a method that divides a path. Therefore, our approach can be implemented on an RDBMS easily, even if long paths exist in some logistic applications and we cannot store a large number in one attribute. Therefore, our approach is efficient and practical. Its application is not limited to supply chain management nor logistic applications. It will be useful for many applications that need flow-based processing. However, since we have not considered a distributed environment, all RFID data should be collected in the central server. Data collection in the central server may be difficult in some environments. Therefore, in the future, we plan to adapt and develop our approach in a distributed environment.

ACKNOWLEDGMENTS

The authors would like to thank the associate editor who handled this paper nicely, and anonymous reviewers for their helpful comments. This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MEST) (No. 2010-0000863).

REFERENCES

- [1] Chinese Remainder Theorem. http://en.wikipedia.org/wiki/Chinese_remainder_theorem, 2010.
- [2] Epc Global. <http://www.epcglobalinc.org/home>, 2010.
- [3] Xpath. <http://www.w3.org/TR/xpath>, 2010.
- [4] R. Agrawal, A. Cheung, K. Kailing, and S. Schönauer, "Towards Traceability across Sovereign, Distributed RFID Databases," *Proc. 10th Int'l Database Eng. and Applications Symp. (IDEAS)*, 2006.
- [5] Y. Bai, F. Wang, and P. Liu, "Efficiently Filtering RFID Data Streams," *Proc. First Int'l Very Large Data Bases (VLDB) Workshop Clean Databases*, 2006.
- [6] Y. Bai, F. Wang, P. Liu, C. Zaniolo, and S. Liu, "RFID Data Processing with a Data Stream Query Language," *Proc. 23rd IEEE Int'l Conf. Data Eng. (ICDE)*, 2007.
- [7] C. Ban, B. Hong, and D. Kim, "Time Parameterized Interval R-Tree for Tracing Tags in RFID Systems," *Proc. 16th Int'l Conf. Database and Expert Systems Applications (DEXA)*, pp. 503-513, 2005.
- [8] C. Bornhövd, T. Lin, S. Haller, and J. Schaper, "Integrating Automatic Data Acquisition with Business Processes - Experiences with SAP's Auto-ID Infrastructure," *Proc. 30th Int'l Conf. Very Large Data Bases (VLDB)*, pp. 1182-1188, 2004.
- [9] D.J. Bowersox and D.J. Closs, *Logistical Management*. McGraw-Hill, 1996.
- [10] P. Buneman, M. Grohe, and C. Koch, "Path Queries on Compressed XML," *Proc. 29th Int'l Conf. Very Large Data Bases (VLDB)*, pp. 141-152, 2003.
- [11] S.S. Chawathe, V. Krishnamurthy, S. Ramachandran, and S. Sarma, "Managing RFID Data," *Proc. 30th Int'l Conf. Very Large Data Bases (VLDB)*, pp. 1189-1195, 2004.
- [12] L. Forlizzi, R.H. Güting, E. Nardelli, and M. Schneider, "A Data Model and Data Structures for Moving Objects Databases," *Proc. ACM SIGMOD*, pp. 319-330, 2000.
- [13] H. Gonzalez, J. Han, and X. Li, "FlowCube: Constructing RFID FlowCubes for Multi-Dimensional Analysis of Commodity Flows," *Proc. 32nd Int'l Conf. Very Large Data Bases (VLDB)*, pp. 834-845, 2006.
- [14] H. Gonzalez, J. Han, X. Li, and D. Klabjan, "Warehousing and Analyzing Massive RFID Data Sets," *Proc. 22nd Int'l Conf. Data Eng. (ICDE)*, 2006.
- [15] R.H. Güting, M.H. Böhlen, M. Erwig, C.S. Jensen, N.A. Lorentzos, M. Schneider, and M. Vazirgiannis, "A Foundation for Representing and Querying Moving Objects," *ACM Trans. Database Systems*, vol. 25, no. 1, pp. 1-42, 2000.
- [16] H. Haidarian-Shahri, G. Namata, S. Navlakha, A. Deshpande, and N. Roussopoulos, "A Graph-Based Approach to Vehicle Tracking in Traffic Camera Video Streams," *Proc. Workshop Data Management for Sensor Networks, in Conjunction with VLDB (DMSN)*, pp. 19-24, 2007.
- [17] J.E. Hoag and C.W. Thompson, "Architecting RFID Middleware," *IEEE Internet Computing*, vol. 10, no. 5, pp. 88-92, Sept./Oct. 2006.
- [18] Y. Hu, S. Sundara, T. Chorma, and J. Srinivasan, "Supporting RFID-Based Item Tracking Applications in Oracle DBMS Using a Bitmap Datatype," *Proc. 31st Int'l Conf. Very Large Data Bases (VLDB)*, pp. 1140-1151, 2005.
- [19] S.R. Jeffery, M.N. Garofalakis, and M.J. Franklin, "Adaptive Cleaning for RFID Data Streams," *Proc. 32nd Int'l Conf. Very Large Data Bases (VLDB)*, pp. 163-174, 2006.
- [20] N.K. Kanhere and S.T. Birchfield, "Real-Time Incremental Segmentation and Tracking of Vehicles at Low Camera Angles Using Stable Features," *IEEE Trans. Intelligent Transportation Systems*, vol. 9, no. 1, pp. 148-160, Mar. 2008.
- [21] Z. Kim and J. Malik, "Fast Vehicle Detection with Probabilistic Feature Grouping and Its Application to Vehicle Tracking," *Proc. Ninth IEEE Int'l Conf. Computer Vision (ICCV)*, pp. 524-531, 2003.
- [22] C.-H. Lee and C.-W. Chung, "Efficient Storage Scheme and Query Processing for Supply Chain Management Using RFID," *Proc. ACM SIGMOD*, 2008.
- [23] Q. Li and B. Moon, "Indexing and Querying XML Data for Regular Path Expressions," *Proc. 27th Int'l Conf. Very Large Data Bases (VLDB)*, 2001.
- [24] W. Lian, N. Mamoulis, D.W.-L. Cheung, and S.-M. Yiu, "Indexing Useful Structural Patterns for XML Query Processing," *IEEE Trans. Knowledge and Data Eng.*, vol. 17, no. 7, pp. 997-1009, July 2005.
- [25] J.-K. Min, M.-J. Park, and C.-W. Chung, "XPRESS: A Queriable Compression for XML Data," *Proc. ACM SIGMOD*, 2003.

- [26] M.F. Mokbel, X. Xiong, and W.G. Aref, "SINA: Scalable Incremental Processing of Continuous Queries in Spatio-Temporal Databases," *Proc. ACM SIGMOD*, pp. 623-634, 2004.
- [27] D. Pfoser, C.S. Jensen, and Y. Theodoridis, "Novel Approaches to the Indexing of Moving Object Trajectories," *Proc. 26th Int'l Conf. Very Large Data Bases (VLDB)*, pp. 395-406, 2000.
- [28] K. Porkaew, I. Lazaridis, and S. Mehrotra, "Querying Mobile Objects in Spatio-Temporal Databases," *Proc. Seventh Int'l Symp. Advances in Spatial and Temporal Databases (SSTD)*, pp. 59-78, 2001.
- [29] J. Rao, S. Doraiswamy, H. Thakkar, and L.S. Colby, "A Deferred Cleansing Method for RFID Data Analytics," *Proc. 32nd Int'l Conf. Very Large Data Bases (VLDB)*, pp. 175-186, 2006.
- [30] P. Rao and B. Moon, "Prix: Indexing and Querying XML Using Pruffer Sequences," *Proc. 20th Int'l Conf. Data Eng. (ICDE)*, 2004.
- [31] K.H. Rosen, *Discrete Mathematics and Its Applications*. McGraw-Hill, 2003.
- [32] Y. Tao and D. Papadias, "MV3R-Tree: A Spatio-Temporal Access Method for Timestamp and Interval Queries," *Proc. 27th Int'l Conf. Very Large Data Bases (VLDB)*, pp. 431-440, 2001.
- [33] I. Tatarinov, S. Viglas, K.S. Beyer, J. Shanmugasundaram, E.J. Shekita, and C. Zhang, "Storing and Querying Ordered XML Using a Relational Database System," *Proc. ACM SIGMOD*, 2002.
- [34] P.M. Tolani and J.R. Haritsa, "XGRIND: A Query-Friendly XML Compressor," *Proc. 18th Int'l Conf. Data Eng. (ICDE)*, pp. 225-234, 2002.
- [35] F. Wang and P. Liu, "Temporal Management of RFID Data," *Proc. 31st Int'l Conf. Very Large Data Bases (VLDB)*, pp. 1128-1139, 2005.
- [36] F. Wang, S. Liu, P. Liu, and Y. Bai, "Bridging Physical and Virtual Worlds: Complex Event Processing for RFID Data Streams," *Proc. 10th Int'l Conf. Extending Database Technology (EDBT)*, 2006.
- [37] H. Wang, S. Park, W. Fan, and P.S. Yu, "ViST: A Dynamic Index Method for Querying XML Data by Tree Structures," *Proc. ACM SIGMOD*, 2003.
- [38] W. Wu, W. Guo, and K.-L. Tan, "Distributed Processing of Moving K-Nearest-Neighbor Query on Moving Objects," *Proc. 23rd IEEE Int'l Conf. Data Eng. (ICDE)*, pp. 1116-1125, 2007.
- [39] X. Wu, M.-L. Lee, and W. Hsu, "A Prime Number Labeling Scheme for Dynamic Ordered XML Trees," *Proc. 20th Int'l Conf. Data Eng. (ICDE)*, 2004.
- [40] C. Zhang, J.F. Naughton, D.J. DeWitt, Q. Luo, and G.M. Lohman, "On Supporting Containment Queries in Relational Database Management Systems," *Proc. ACM SIGMOD*, 2001.



Chun-Hee Lee received the PhD degree in computer science from the Korea Advanced Institute of Science and Technology (KAIST), Korea, in 2010. His research interests include sensor network, stream data management, and graph databases.



Chin-Wan Chung received the BS degree in electrical engineering from Seoul National University, Korea, in 1973, and the PhD degree in computer engineering from the University of Michigan, Ann Arbor, in 1983. From 1983 to 1993, he was a senior research scientist and a staff research scientist in the Computer Science Department at the General Motors Research Laboratories (GMR). While at GMR, he developed Dataplex, a heterogeneous distributed

database management system integrating different types of databases. Since 1993, he has been a professor in the Department of Computer Science at the Korea Advanced Institute of Science and Technology (KAIST), Korea. At KAIST, he developed a full-scale object-oriented spatial database management system called OMEGA, which supports ODMG standards. His current major project is about mobile social networks in Web 3.0. He was in the program committees of major database conferences including ACM SIGMOD, VLDB, and IEEE ICDE. He was an associate editor of *ACM TOIT*, and is an associate editor of *WWW Journal*. His current research interests include the semantic Web, mobile Web, sensor network and stream data management, and multimedia databases. More information is available at <http://islab.kaist.ac.kr/chungcw>.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.