



Volume 5, Issue 6

August 2007

ISSN 1570-8705

Ad Hoc Networks

Special Issue

(1) Wireless Mesh Networks

Guest Editors:

X. Wang, E. Knightly, M. Conti and
A. Ephremides

(2) Wireless Sensor Networks

Guest Editors:

E. Ekici and L. Kleinrock

This article was originally published in a journal published by Elsevier, and the attached copy is provided by Elsevier for the author's benefit and for the benefit of the author's institution, for non-commercial research and educational use including without limitation use in instruction at your institution, sending it to specific colleagues that you know, and providing a copy to your institution's administrator.

All other uses, reproduction and distribution, including without limitation commercial reprints, selling or licensing copies or access, or posting on open internet sites, your personal or institution's website or repository, are prohibited. For exceptions, permission may be sought for such use through Elsevier's permissions site at:

<http://www.elsevier.com/locate/permissionusematerial>

Communication-efficient implementation of join in sensor networks

Himanshu Gupta^{*}, Vishal Chowdhary

Department of Computer Science, Stony Brook, NY 11794, United States

Received 10 February 2007; accepted 15 February 2007

Available online 24 February 2007

Abstract

A sensor network is a multi-hop wireless network of sensor nodes cooperatively solving a sensing task. Each sensor node generates data items that are readings obtained from one or more sensors on the node. This makes a sensor network similar to a distributed database system. While this view is somewhat traditional, efficient execution of database (SQL) queries in sensor network remains a challenge, due to the unique characteristics of such networks such as limited memory and battery energy on individual nodes, multi-hop communication, unreliable infrastructure, and dynamic topology. Since the nodes are battery powered, the sensor network relies on energy-efficiency (and hence, communication efficiency) for a longer lifetime of the network.

In this article, we have addressed the problem of communication-efficient implementation of the SQL **join** operator in sensor networks. In particular, we design an optimal algorithm for implementation of a join operation in dense sensor networks that provably incurs minimum communication cost under some reasonable assumptions. Based on the optimal algorithm, we design a suboptimal heuristic that empirically delivers a near-optimal join implementation strategy and runs much faster than the optimal algorithm. Through extensive simulations on randomly generated sensor networks, we show that our techniques achieve significant energy savings compared to other simple approaches.

© 2007 Elsevier B.V. All rights reserved.

Keywords: Sensor network databases; Join implementation; Communication efficiency

1. Introduction

A sensor network consists of sensor nodes with a short-range radio and on-board processing capability forming a multi-hop network of an irregular topology. Each sensor node can sense certain phys-

ical phenomena like light, temperature, or vibration. There are many exciting applications of such sensor networks, including monitoring and surveillance systems in both military and civilian contexts, building smart environments and infrastructures such as intelligent transportation systems and smart homes.

Each sensor node typically generates a stream of data items that are readings obtained from one or more sensing devices on the node. This motivates visualizing sensor networks as distributed database systems[8,15,18] and the data present in a sensor

^{*} Corresponding author. Tel.: +1 631 632 8446; fax: +1 631 632 8334.

E-mail addresses hgupta@cs.suny.edu (H. Gupta), vishalc@microsoft.com (V. Chowdhary).

network as relational data streams. Like a database, the sensor network is queried to gather the sensed data tuples. Database queries in SQL are a very general representation of queries over data, and because of the enormous amount of data present in a typical sensor network, efficient implementation of database queries is of great significance.

The main performance criterion for distributed implementations of queries in sensor network is the total communication cost incurred, since each sensor node has limited battery power and message communication nodes is the main consumer of battery energy. Thus, distributed implementation of queries must minimize the communication cost incurred. In particular, we are interested in in-network implementation strategies since a centralized strategy of transmitting all sensor data to a central server for further computation would incur prohibitive communication costs.

In this article, we focus on designing efficient distributed implementations for the join operation in sensor networks. The join operator is essentially a cartesian product of the operand tables followed by a predicate selection. The motivation for the join operation in sensor networks comes from one of the most prominent sensor network applications viz., event detection, wherein complex events can be defined as joins over data streams [2,11]. We propose a novel path-join algorithm, which computes the join result by first distributing one of the operand tables along a predetermined path of sensors. Using path-join algorithm as the basic step, we design an optimal algorithm for a join operation that provably incurs minimum communication cost in dense sensor networks under some reasonable assumptions of communication cost and computation model. We also design a much faster suboptimal heuristic that empirically performs very close to the optimal algorithm, and results in significant savings over the naive approaches.

1.1. Paper organization

The rest of the paper is organized as follows. We start with modeling the sensor network as a database and motivating implementation of the join operation in the sensor network. In Section 3, we present various algorithms for in-network implementation of the join operator for static (non-streaming tables). In Section 4, we generalize our techniques to handle streaming tables and discuss relaxation of other assumptions. We present our

experiment results in Section 5. Related work is discussed in Section 6, and concluding remarks presented in Section 7.

2. Sensor network databases

A sensor network consists of a large number of sensors distributed randomly in a geographical region. Each sensor has limited on-board processing capability and is equipped with sensing devices. We assume that each sensor node is aware of its geographic location (obtained using GPS or other localization techniques [5]). A sensor node also has a radio which is used to communicate directly with some of the sensors around it. Two sensor nodes can communicate with each other if and only if the distance between them is less than the transmission radius. We assume that each sensor node in the sensor network has a limited storage capacity of m units. As mentioned above, each sensor node has limited battery energy, which must be conserved for prolonged unattended operation. Thus, we have focused on minimization of communication cost (hence, energy cost) as the key performance criteria of the join implementation strategies.

2.1. Modeling the sensor network as a database

In a sensor network, the data generated by the sensor nodes is simply the readings obtained from the sensing devices on the node. The data records produced by a group of sensor nodes with similar capabilities and responsibility will have similar format and semantics, and thus, can be modeled as rows of the same relational table. More specifically, due to the continuous generation of data tuples in the sensor network, the sensor network data is best modeled as data streams [3]. The above motivates visualizing sensor networks as distributed database systems [8,15,18] of streaming tables. In a sensor network, a data stream may be partitioned horizontally across (or generated by) a set of sensors in the network. Each data stream has a corresponding generating region which could very well be the entire network region. Due to the spatial and real-time nature of the data generated, a tuple usually has `timeStamp` and `nodeLocation` as attributes, and the sensor node that generates a particular tuple is referred as its source node. Like traditional database systems, the sensor network database can also be queried to access and manip-

ulate the data tables, and SQL with some extensions can be used as a query language for sensor networks.

2.1.1. Database queries

A database query is composed of one or more database operators. The core database operators are viz. selection (selecting tuples based on a predicate), projection (selecting given attributes of a table), join (cartesian product followed by selection), grouping (partitioning a table based on a set of attribute values), aggregation (aggregating attributes for each group), outerjoins (join plus the unmatched tuples padded with NULLs), duplicate elimination, union, difference, and intersection. Union, difference, and intersection have same semantics as the corresponding set operators.

The focus of this article is communication-efficient in-network implementation of the join operator. The join operator is used to correlate data from multiple tables and is essentially a cartesian product of the operand tables followed by a selection. As selection and projection are unary operators and operate on each tuple independently, they could be implemented by computing the operation locally followed by efficiently routing to the query source. Union operation can be reduced to duplicate elimination, and the difference and intersection operations can be reduced to the join operation. Implementation of other database operators (aggregation, duplicate elimination, and outerjoins) is challenging and is part of our future work.

2.1.2. In-network implementation of SQL queries

A plausible implementation of a sensor network database query engine could be to have an external database system handle all the queries over the network. In such a realization, all the data from each sensor node in the network is sent to the external system that handles the execution of queries completely. Such an implementation would incur very high communication costs and congestion-related bottlenecks. Thus, prior research has proposed query engines that would execute the queries within the network with little external help. In particular, [9] shows that in-network implementation of database queries is fundamental to achieving energy-efficient communication in sensor networks. Moreover, due to the very limited processing memory available on a sensor nodes, it will be impossible to compute

the join locally on any particular node, especially for large tables.

2.1.3. Querying and cost model in sensor networks

A query in a sensor network is initiated at a node called query source and the result of the query is required to be routed back to the query source for storage and/or consumption. A stream database table may be generated by a set of sensor nodes in a closed geographical region. The optimization algorithms, proposed in this article, to determine how to implement the join operation efficiently, are run at the query source. As typical sensor network queries are long running, the query source can gather all the catalogue information needed (estimated sizes and locations of the operand relations, join selectivity factor to estimate the size of the join result, density of the network) by initially sampling the operand tables. As mentioned before, we concentrate on implementations that minimize communication cost. We define the total communication cost incurred as the total data transfer between neighboring sensor nodes.

Our algorithms target the general long-running queries in the sensor network. Given a query source Q and regions R and S where a join has to be taken. Initially all the tuples of the participating tables are routed to the query source Q , which collects catalog information and estimates parameters such as locations of the region R and S , sizes of R , S and $R \cap S$, and join selectivity factor f . Using the optimal algorithm, the query source Q calculates the optimal region P where the join should be executed in the sensor network.

2.1.4. Join in sensor networks

The SQL join operator is used to correlate data from multiple tables, and can be defined as a selection (join) predicate over the cross-product of a pair of tables; a join of R and S tables is denoted as $R \Join S$. One of the most popular applications of sensor networks is event detection which motivates the body of our work. An event indicates a point in time of interest based on certain conditions over the generated sensor data. For certain applications, events may simply depend on the local value of a particular sensor reading. Higher-level events or complex events may be specified using composition operators over the primitive events. In particular, the complex events may be represented as a join of multiple data streams, involving spatial and temporal constraints and correlations.

3. In-network implementation of join

In this section, we first develop communication-efficient algorithms for implementation of a join operation over static (non-streaming) database tables stored in some sensor network region. As data in sensor network is better represented as data stream tables, we will generalize our techniques for stream database tables in the next section.

Consider a join operation, initiated by a query source node Q , involving two static (non-streaming) tables R and S distributed horizontally across some geographical regions R and S in the network. We assume that the geographic regions are disjoint and small relative to the distances between the query source and the operand table regions. We later discuss generalizing our algorithms for general query source locations and operand regions. If we do not make any assumptions about the join predicates involved, each data tuple of table R should be paired with every tuple of S and checked for the join condition. The joined tuple is then routed (if it passes the join selection condition) to the query source Q where all the tuples are accumulated or consumed. Given that each sensor node has limited memory resources, we need to find out appropriate regions in the network that would take the responsibility of computing the join. In particular, we may need to store and process the relations at some intermediate location before routing the result to the query source.

A simple nested-loop implementation of a join used in traditional databases is to generate the cross-product (all pairs of tuples), and then extract those pairs that satisfy the selection predicate of the join. More involved implementations of a join operator widely used in database systems are merge-sort and hash-join. These classical methods are unsuitable for direct implementation in sensor networks due to the limited memory resources at each node in the network. Moreover, the traditional join algorithms focus on minimizing computation cost, while in sensor networks the primary performance criteria is communication cost. Below, we discuss various techniques for efficient implementation of the join operation in sensor networks.

Naive approach A simple way to compute $R \bowtie S$ could be to route the tuples of S from their original location S to the region R , broadcast the S -tuples in the region R , compute the join within the region R , and then route the joined tuples to the query source

Q . The breakup of the total communication cost incurred is as follows:

- (1) Cost incurred in routing the table S to the region R .
- (2) Cost incurred in broadcasting the table S throughout R .
- (3) Cost incurred in routing the result (from R) to the query source Q .

Note that in the above approach the roles of the tables R and S can be interchanged.¹

Centroid approach. Now, we consider another approach where the region responsible for computing the join operation is a circular region around some point C in the sensor network. Let $|R|$ denote the size of the table R , m denote the memory of each sensor node, and let P_c be the smallest circular region around C such that the region P_c has more than $|R|/m$ sensor nodes to store the table R . First we route and distribute the tuples of table R in the region P_c , and then route and broadcast the tuples of table S in the region P_c . After computing the join operation in the region P_c , we route the resulting tuples of $R \bowtie S$ to the query source Q . The communication cost incurred consists of the following components. (i) Cost incurred in routing the tables to C . (ii) Cost incurred in distributing R and broadcasting S in the region P_c around C . (iii) Cost incurred in routing the result $R \bowtie S$ to the query source Q . Since the second component of the cost is independent of the choice of C , it is easy to see that the communication cost incurred in the above approach is minimized when the point C is the weighted centroid of the triangle formed by R , S , and Q (i.e., the point that minimizes the sum of the weighted distances from the three points). Here, the choice of the centroid point C is weighted by the sizes of R , S , and $R \bowtie S$.

3.1. Path-join algorithm

In the above two paragraphs, we described a couple of simple approaches to compute the join operation in a sensor network. However, in order to minimize communication cost, we may need to perform the join operation in a region having a non-trivial shape. In this subsection, we present a

¹ The other simple approach of computing the join at a region around Q is subsumed by the Centroid Approach discussed next.

novel path-join approach of performing a join operation. In the next subsection, we will extend the path-join approach to devise an optimal join algorithm that incurs minimum communication cost in dense sensor networks.

The path-join implementation of the join operation works as follows. First, all the tuples of R are distributed uniformly along an appropriately chosen path P containing $|R|/m$ sensor nodes, where $|R|$ is the size of table R and m is the memory size of each sensor node. Then, every tuple t is routed to the path P and passed through all the sensors along P to perform the join. The resulting joined tuples computed at each sensor along the path P are then routed to the query source Q . See Fig. 1. The location of the path P is chosen to minimize the total communication cost incurred. We estimate the total communication cost incurred in terms of a notion of sensor length, defined below.

Definition 1 (Sensor length $d(X; y)$ and notation $|X|$). The sensor length between a region X and a point y in a sensor network plane is denoted as $d(X; y)$ and is defined as the average weighted distance, in terms of number of hops (i.e., intermediate nodes), between the region X and the point y . Here, the distance between a point $x \in X$ and y is weighted by the amount of data residing at x .

For a region X in the sensor networks, the notation $|X|$ denotes the number of sensors in the region $|X|$. Note that for a relational table R , we use $|R|$ to denote the size of the table R .

Let $|R|$, $|S|$, and $|R \bowtie S|$ be the respective sizes of the tables R , S , and the joined result $R \bowtie S$. Let Q be the query source C_0 be an end of the path P that is closer to R and/or S , and $|P|$ be the number of sensors on the path P . Note that by choice of P , $|P| \cdot \frac{1}{m} |R| = m$. Let us assume that both R and S start their broadcast and distribution phases from the same point C_0 . The total communication cost incurred in the path-join algorithm consists of: cost

of routing R to C_0 , cost of routing S to C_0 , cost of distributing the table R along the path P , cost of broadcasting the table S along the path P , and finally, the cost of routing the joined tuples from P to Q . If we assume that the resulting joined tuples are uniformly distributed along the path P , then the total communication cost incurred is

$$|R|d(C_0; R) + |S|d(C_0; S) + |P| \cdot \frac{|R|}{m} d(C_0; P) + |P| \cdot |S| d(P; Q)$$

Note that the distribution and broadcast cost of R and S respectively is independent of the location of P . In some cases, the path-join algorithm may not be optimal, i.e., may incur more than the minimum communication cost possible.

3.2. Optimal join algorithm

In this section, we present an algorithm that uses path-join as a basic component, and constructs a region for computing the join operation using optimal communication cost. We assume that the sensor network is sufficiently dense that we can find a sensor node at any point in the region. To formally prove the claim of optimality, we need to restrict ourselves to a class of join algorithms called **Distribute-Broadcast Join Algorithms** (defined below). In effect, our claim of optimality states that the proposed join algorithm incurs less communication cost than any distribute-broadcast join algorithm.

Definition 2 (Distribute-broadcast join algorithms). A join algorithm to compute $R \bowtie S$ in a sensor network is a distribute-broadcast join algorithm if the join is processed by "first uniformly distributing the table R in some region P (other than the region R storing R)² of the sensor network followed by broadcasting the relation S within the region P to compute the join. The joined tuples are then routed from each sensor in the region P to the query source.

As before, consider a query source Q , and regions R and S that store the static operand tables R and S in a sensor network. The key challenge in designing an optimal algorithm for implementation of a join operation is to select a region P for processing the join in such a way that the total communication cost is minimized. Note that in general, P

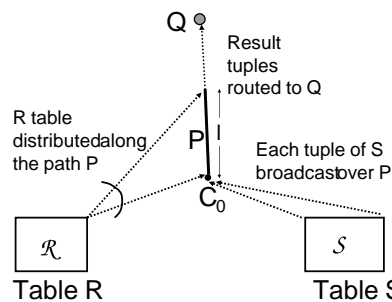


Fig. 1. Path-join implementation. Here, $|P| \cdot \frac{1}{m} |R| = m$.

² Else, the algorithm will be identical to one of the naive approaches.

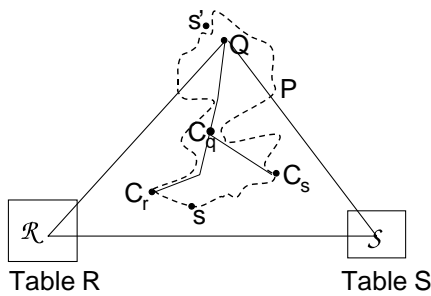


Fig. 3. An arbitrary join-region P containing Q , and the shortest paths $\tilde{d}_{C_r;Q}$ and $\tilde{d}_{C_s;Q}$ in P . Here, s is the point in $P \cap \tilde{d}_{C_r;Q} \cap \tilde{d}_{C_s;Q}$ that is farthest from Q , and s^0 is the point not in P that is closest to Q . If P is not of the form Fig. 2a, then $d_{P;Q}$ can be reduced (without changing C_r, C_s , or $|P|$) by replacing s by s^0 , and thus, reducing the total cost $T_{C_r;C_s;P}$.

of the total cost $T_{C_r;C_s;P}$ since we keep $|P|$, C_r , and C_s fixed.

Reducing $d_{P;Q}$ when $Q \in P$. Let us consider the case when Q is in P , but the region P is not of the form Fig. 2a. Let C_r and C_s be the points in P to where R and S are routed for distribution and broadcast respectively in the region P . See Fig. 3. Consider paths (not necessarily disjoint) $\tilde{d}_{C_r;Q}$ and $\tilde{d}_{C_s;Q}$ contained in P that connect C_r and C_s respectively to Q using minimum number of sensor nodes. Since P is connected, such paths exist. Consider a point s in P such that s is neither on path $\tilde{d}_{C_r;Q}$ nor on path $\tilde{d}_{C_s;Q}$ and is farthest away from Q . Such a point s must exist, else P would be comprised entirely of paths $\tilde{d}_{C_r;Q}$ and $\tilde{d}_{C_s;Q}$ and hence, of the form Fig. 2a.⁴

Now, consider a point $s^0 \notin P$ that is closest to Q . See Fig. 3. If Q is closer to s than s^0 , i.e., if $d_{s;Q} < d_{s^0;Q}$, then P is just comprised of the paths $\tilde{d}_{C_r;Q}$, $\tilde{d}_{C_s;Q}$, and a fully packed circular region around Q , and thus, of the form Fig. 2a. Since we assumed to the contrary Q must be closer to s^0 than s , i.e., $d_{s^0;Q} < d_{s;Q}$. Now, in such a case, $d_{P;Q}$ can be reduced as follows. Since $d_{P;Q} = \frac{1}{4}|P| + \int_{P \cap \tilde{d}_{P;Q}} d_{P;Q}$ (since R is uniformly distributed in P), the value $d_{P;Q}$ can be reduced by changing P to $P \setminus \{s\} \cup \{s^0\}$, i.e., replacing s by s^0 . Note that such a point s^0 will be directly connected to P , and hence, addition of s^0 maintains the connectivity of P . Moreover, since s is neither in

$\tilde{d}_{C_r;Q}$ nor in $\tilde{d}_{C_s;Q}$ and is farthest such point from Q , removal of s from P maintains the connectivity of P . Finally, the above replacement keeps C_r, C_s , and $|P|$ fixed.

Final arguments. Thus, we can reduce $d_{P;Q}$ while keeping C_r, C_s , and $|P|$ fixed, and thus, reduce the total cost $T_{C_r;C_s;P}$ when $Q \in P$ and P is not of the form Fig. 2a. Thus, by contradiction, the optimal join-region P must be of the form depicted in Fig. 2a if $Q \in P$. Using similar arguments, we can show that if $Q \notin P$, the optimal join-region must be of the shape depicted in Fig. 2b. \square

Note that the assumptions made (viz. restricted class of algorithms, distributing and broadcasting in a linear fashion) in proving the above theorem do not restrict the applicability of our developed techniques. The assumptions were made solely to prove optimality, and more importantly, to develop an algorithm that could form the basis of a communication-efficient implementation of the join operation in general sensor networks without any restrictions on the communication/computation model.

Note that the above theorem only restricts the shape of an optimal join-region; there are still an infinite number of possible join-regions of shapes depicted in Fig. 2. Thus, we now further restrict the shape of an optimal join-region by characterizing the equations of the paths P_r and P_s that connect C_r and C_s respectively to C_q .

3.2.2. Optimizing paths P_r and P_s in the join-region

Consider an optimal join-region P that implements a join operation using minimum communication cost. By Theorem 1, we know that the region P is of the shape depicted in Fig. 2a or b. As derived in Eq. (1), the total communication cost $T_{C_r;C_s;P}$ incurred in processing of a join using the region P is $\int_R \int_S d_{P;Q} + \int_{C_r} \int_P d_{P;Q} + \int_{C_s} \int_P d_{P;Q} = 2 \int |P| + \int_{C_r} \int_P d_{P;Q} + \int_{C_s} \int_P d_{P;Q}$. Let $P \subset P_r \cup P_s$, i.e., the region P without the paths P_r and P_s . Since the result $\int_R \int_S$ is uniformly spread along the entire region P , we have

$$d_{P;Q} = \frac{1}{|P|} \int_P d_{P;Q} \leq \int_{P_r} d_{P;Q} + \int_{P_s} d_{P;Q}$$

For a given $|P|$ and a given set of points C_r, C_s , and C_q , the total communication cost T is minimized when the path P_r is constructed such that $\int_{P_r} d_{P;Q}$ is minimized. Otherwise, we could reconstruct P_r with a smaller $\int_{P_r} d_{P;Q}$ and

⁴ Note that since paths $\tilde{d}_{C_r;Q}$ and $\tilde{d}_{C_s;Q}$ are shortest in P , they intersect at only one point C_q and have the same subpaths $\tilde{d}_{C_q;Q}$.

remove/add sensor nodes from the end of the region P^0 to maintain $|P_j|$. Removal of sensor nodes from P^0 will always reduce T , and it can be shown that addition of sensor nodes to the end of the region P^0 will not increase the cost more than the reduction achieved by optimizing P_r . Similarly, the path P_s could be optimized independently.

We now derive the equation of the path P_r that minimizes $|P_r| + d(P_r, Q)$ for a given C_r and C_q . Consider an arbitrary point $R(x, y)$ along the optimal path P_r . The length of an infinitesimally small segment of the path P_r beginning at $R(x, y)$ is

$\sqrt{dx^2 + dy^2}$, and the average distance of this segment from Q is $\sqrt{x^2 + y^2}$, if the coordinates of Q are $(0, 0)$. Sum of all these distances over the path P_r is: $F = \int_{C_r}^{C_q} \sqrt{x^2 + y^2} \sqrt{dx^2 + dy^2}$. To get the equation for the path P_r , we would need to determine the extremals of the above function F . Using the technique of calculus of variations [7], we can show that the extremal values of F satisfy the Euler-Lagrange differential equation. The equation of the path P_r can thus be computed as (we omit the details): $b + x^2 \cos a + 2xy \sin a - y^2 \cos a$ where the constants a and b are evaluated by substituting for coordinates of C_r and C_q in the equation.

3.2.3. Computing communication cost

Given $|P_j|$ and the three points C_r , C_s , and C_q , we now derive the total communication cost $T_{opt}(C_r, C_s, C_q; |P_j|)$ incurred by using the optimal join-region of size $|P_j|$ constructed over C_r , C_s and C_q . We will use the formulation of $T_{opt}(C_r, C_s, C_q; |P_j|)$ to design an optimal algorithm by consider all possible combinations of values of $|P_j|$, C_r , C_s and C_q and picking the quartet that results in minimum $T_{opt}(C_r, C_s, C_q; |P_j|)$.

Given $|P_j|$ and points C_r, C_s, C_q , let P_r and P_s be the paths as obtained in the previous paragraph. Let

$$l_y = |P_r| + |P_s| + |C_q Q|;$$

If $l_y > |P_j|$, then the optimal join-region P cannot contain the point Q , and hence, by Theorem 1, the region P is comprised of the optimized paths P_r , P_s , and the line segment $\overline{C_q C_{q2}}$, where $C_{q2} \in \overline{C_q Q}$ is such that $|C_q C_{q2}| = |P_j| - |P_r| - |P_s|$. See Fig. 2b. For the case when $l_y \leq |P_j|$, the $l_y = |P_j|$ frac-

tion of the join is processed on the curve P_r, P_s , and the line segment $\overline{C_q Q}$, while the remaining fraction of the join is processed on a circular region P_o of appropriate radius around Q . See Fig. 2a. From Theorem 1, the above choice of P minimizes the value $d(P, Q)$ for a given combination of C_r, C_s, C_q ; and $|P_j|$. Thus, we have

$$P = P_r \cup P_s \cup \overline{C_q C_{q2}} \quad \text{if } l_y > |P_j|; \quad (2)$$

$$P = P_r \cup P_s \cup \overline{C_q Q} \cup P_o \quad \text{if } l_y \leq |P_j|; \quad (3)$$

As mentioned before, the point C_{q2} is such that $|C_q C_{q2}| = |P_j| - |P_r| - |P_s|$, and P_o is a circular region of sufficient radius around Q such that $|P_o| = |P_j| - |C_q C_{q2}| - |P_r| - |P_s|$. For a given quartet of values $(C_r, C_s, C_q; |P_j|)$, let $T_{opt}(C_r, C_s, C_q; |P_j|)$ denote the total communication cost incurred when the join-region P is optimally constructed as suggested by Eqs (2) and (3). In other words, $T_{opt}(C_r, C_s, C_q; |P_j|)$ is equal to $|R_j| + |C_r| + |S_j| + |C_s| + |P_j|$, where P is the optimally constructed join-region as suggested by Eqs (2) and (3).

3.2.4. Optimal join algorithm

Based on the above discussion, we construct an optimal join-region to compute a join operation for tables R and S and the query source Q , by considering all possible triples of points C_r, C_s , and C_q in the sensor network and values of $|P_j|$, and pick the quartet $(C_r, C_s, C_q; |P_j|)$ that minimizes the value $T_{opt}(C_r, C_s, C_q; |P_j|)$. For such an optimal quartet $(C_r, C_s, C_q; |P_j|)$, we construct the optimal join-region P as suggested by Eqs (2) and (3) in the previous paragraph. If n is the total number of network nodes, then there are at most n^4 combinations of $(C_r, C_s, C_q; |P_j|)$. Thus, the time complexity of the above algorithm which constructs an optimal join-region is $O(n^4)$.

3.2.5. Suboptimal heuristic

The high time complexity of the optimal algorithm described above makes it impractical for large sensor networks. Here, we design a suboptimal heuristic that has a much lower time complexity and performs very well in practice (see Fig. 4). In particular, we reduce the complexity of our designed algorithm from $O(n^4)$ to $O(n^3)$ using the following "ve steps. (i) We choose the minimum value of $|P_j|$, i.e., $|P_j| = |R_j| = m$, where $|R_j|$ is the size of the table R to be distributed and m is the memory at each sensor node. (ii) We look at all possible values for C_r in

⁵ Here, by the end of the region P^0 , we mean either the circular part P_o or the line segment $\overline{C_q C_{q2}}$ depending on the shape.

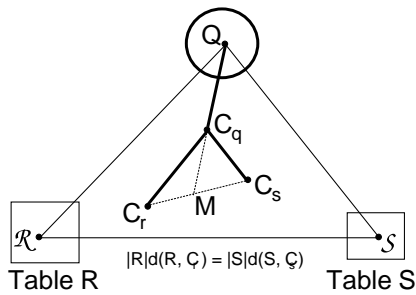


Fig. 4. Suboptimal heuristic for join implementation.

the region. (iii) For each C_r , we stipulate that C_s should be symmetrically located in the region S such that $|R|d(R, C_r) = |S|d(S, C_s)$. Thus, the location of C_s is fixed for a given C_r . (iv) We approximate paths P_r and P_s to be straight line segments $\overline{C_r C_q}$ and $\overline{C_s C_q}$ respectively. (v) We further stipulate that the point C_q should lie on the median of the triangle $MC_r C_s Q$. Thus, for each point as C_r in the sensor network, we determine C_s and search for the best C_q on the median of $MC_r C_s Q$. The above reduces the time complexity to construct a join-region to $O(n^3)$ where n is the network size.

4. Generalizations to stream tables and general sensor networks

In this section, we extend our proposed algorithms to real sensor networks and relax the assumptions made in the previous section. We start with generalizing our technique for stream database tables. Then, we present the overall working of our approach in general sensor networks. Finally, we discuss a few other generalizations.

4.1. Implementation for stream database tables

In the previous section, we discussed implementation of the join operation in a sensor network for static database tables. Since, sensor network data is better represented as stream database tables, we now generalize the algorithms to handle stream database tables. First, we start with presenting our model of stream database tables in sensor networks.

4.1.1. Data streams in sensor networks

As for the case of static tables, a stream database table R corresponding to a data stream in a sensor network is associated with a region R , where each node in R is continually generating tuples for the table R . To deal with the unbounded size of stream

database tables, the tables are usually restricted to a finite set of tuples called the sliding window [1,6,16]. In effect, we expire or archive tuples from the data stream based on some criteria so that the total number of stored tuples does not exceed the bounded window size. We use W_R to denote the sliding window for a stream database table R .

4.1.2. Naive approach for stream tables

In the naive approach, we use the region R (or S) to store the windows W_R and W_S of the stream tables R and S .⁶ Each sensor node in the region R uses $W_R = \alpha |W_R| + \beta |W_S|$ fraction of its local memory to store tuples of W_R , and the remaining fraction of the memory to store tuples of W_S .⁷ We need to store W_S also in the region R to find matches for a newly generated tuple of R . To perform the join operation, each newly generated tuple (of R or S) is broadcast to all the nodes in the region R , and is also stored in some node of R with available memory. Note that the generated data tuples of S need to be first routed from the region S to the region R . The resulting joined tuples are routed from R to the query source Q .

4.1.3. Generalizing other approaches

The other approaches viz. centroid approach, optimal algorithm, and suboptimal heuristic, use a join-region that is separate from the regions R and S . These algorithms are generalized to handle stream database tables as follows. First, the strategy to choose the join-region P remains the same as before for static tables, except for the size of the join-region. For stream database tables, the chosen join-region is used to store W_R as well as W_S , with each sensor node in the join-region using $W_R = \alpha |W_R| + \beta |W_S|$ fraction of its memory to store tuples of W_R , and the rest to store tuples of W_S . We need to store W_S as well in the join-region in order to find matches for the newly generated tuples of R . Now, each newly generated tuple (of R or S) is routed from its source node in R or S to the join-region P , and broadcast to all the nodes in P . The resulting joined tuples are then routed to Q . As part

⁶ If the total memory of the nodes in R is not sufficient to store W_R and W_S , then the region R is expanded to include more sensor nodes.

⁷ An alternate naive strategy could be to store W_R and W_S in R and S respectively, but route each new tuple of R to S and each new tuple of S to R . Such a strategy uses more number of nodes for storages, but incurs more routing communication cost.

of the broadcast process (without incurring any additional communication cost), each generated tuple of R (or S) is also stored at some node in P with available memory.

4.2. Overall implementation in real sensor networks

In this subsection, we consider overall working of our approaches in general sensor networks. We start with discussing the construction of join-region and details of the underlying routing protocols appropriate for our developed techniques.

4.2.1. Join-regions and routing protocols in general networks

Till now, we have assumed “geometric” sensor networks, and looked at the problem of “finding an optimal join-region in a geometric sense. In other words, we assumed that the sensor network is very dense so that we can “find a sensor node at any desirable point in the region. In case of non-geometric (i.e., not sufficiently dense) networks, we define the join-region based on the paths traversed by appropriate routing protocols. In particular, we use GPSR [10] and TBF (trajectory based forwarding [17]) routing protocols to traverse appropriate parts of the intended join-region. More specifically, we use the paths traversed by GPSR protocol as the paths for the line-segment parts of the join-region, i.e., $\overline{C_q Q}$ (or $\overline{C_q C_{q2}}$), and the paths P_r and P_s in the suboptimal heuristic. However, for the curved (non-straight) parts of the join-region (i.e., the paths P_r and P_s in the optimal algorithm), we need to use the TBF technique, which works by forwarding packets to nodes closest to the intended path/trajectory. For reasonably dense sensor networks, the above approach yields a join-region that is very close to the originally intended optimal geometric join-region.

4.2.2. Overall working of our approaches

Recall that the algorithms to construct the join-regions are run at the query source. As typical sensor network queries are long running, the query source can gather all the catalogue information needed (estimated sizes and locations of the operand relations, join selectivity factor, network density) by initially sampling the operand tables. When the query source Q needs to issue a join query, it determines the join-region based on the catalogue information, and passes the constructed join-region (represented by the paths P_r , P_s , and $\overline{C_q C_{q2}}$ (or

$\overline{C_q Q}$ and radius around Q)) to all the nodes in the regions R and S . Each generated tuple of stream R is routed from its source node (in region R) to the node nearest to C_r using GPSR protocol. On reaching C_r , we use GPSR/TBF protocol to route the tuple r through the path P_r to reach the node nearest to C_q , and then use GPSR to route to the node nearest to C_{q2} or Q depending on the join-region. Finally, if needed, the tuple is broadcast in a region around Q of appropriate radius. In addition, during the above traversal, the tuple is joined with tuples of W_s (the sliding window of S) stored locally at each node of the join-region. Also, the tuple r is stored at the “first encountered node with available memory in the join-region.

4.2.2.1. Effect of node failures As described above, our proposed implementations do not use any specific destination nodes for traversing the constructed join-region. That is, even though the join-region is originally represented by certain geographic locations and paths, the actual join-region traversed is based on the paths traversed by GPSR/TBF protocols to nodes nearest to geometric locations. Thus, our overall techniques automatically adapt to node failures just as the underlying routing protocols.

5. Performance evaluation

In this section, we present our simulation results comparing performance of various algorithms designed in this article. In particular, we compare the performance of Naive Approach, Centroid Algorithm, Optimal Algorithm, and Suboptimal Heuristic. Each algorithm is generalized for stream database tables and non-geometric general sensor network. We refer to the generalized algorithms as Naive, Centroid, OptBased and Suboptimal Heuristic respectively. Our simulations demonstrate the effectiveness of our developed techniques. We start with defining join-selectivity factor which is used to characterize the size of the join result.

Definition 3 (Join-selectivity factor). Given instances of relations R and S and a join predicate, the join-selectivity factor is the probability that a random pair of tuples from R and S will satisfy the given join predicate. In other words, the join selectivity factor is the ratio of the size of $R \Join S$ to the size of the cartesian product, i.e., $|R \Join S| = \sigma_j(R \Join S) / |R| \times |S|$.

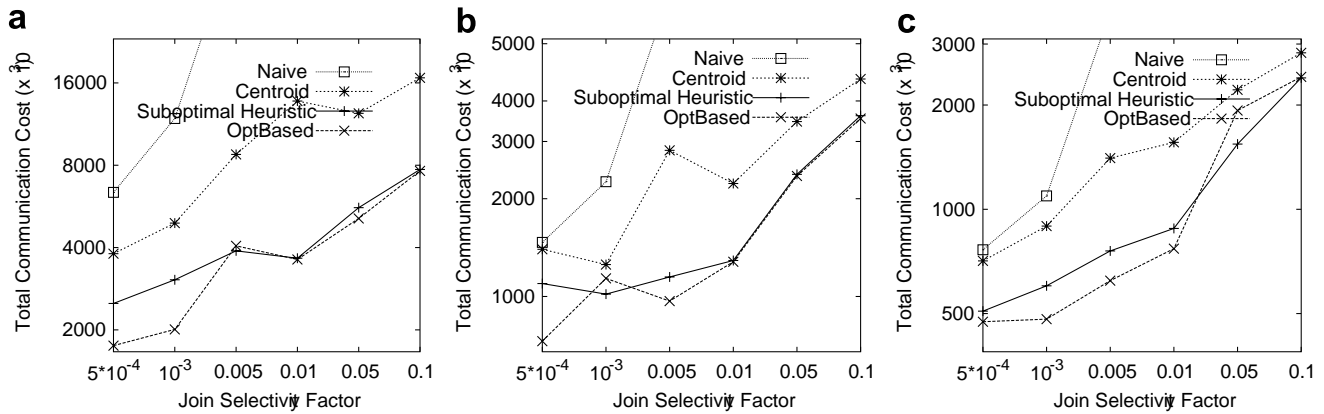


Fig. 5. Performance for a “xed MRSQ with varying join-selectivity factor for three different transmission radii: (a) 0.13 units, (b) 0.15 units, (c) 0.18 units.

5.1. Parameter values and experiments

We generated random sensor networks by randomly placing 10,000 sensors in an area of 1010 units. Each sensor has a uniform transmission radius, and two sensors can communicate with each other if they are located within each other’s transmission radius. For the purposes of comparing the performance of our algorithms, varying the number of sensors is tantamount to varying the transmission radius. Thus, we “x the number of sensors to be 10,000 and measure performance for different transmission radii. Memory size of a sensor node is 300 tuples, and the size of each of the sliding windows W_R and W_S of stream tables R and S is 8000 tuples. For simplicity, we chose uniform data generation rates for R and S streams. In each of the experiments, we measure communication cost incurred in processing 8000 newly generated tuples of R and S each, after the join-region is already “lled with previously generated tuples. We use the GPSR [10] algorithm to route tuples. Catalogue information is gathered for non-Naive approaches by collecting a small sample of data streams at the query source.

We ran three sets of experiments on randomly generated sensor networks. In the “rst set of experiments, we consider a “xed MRSQ and calculate the total communication cost for various transmission radii and join-selectivity factors. Next, we “x the transmission radius and calculate the total communication cost for various join-selectivity factors and various shapes/sizes of the MRSQ. Finally, we plot of performance of various algorithms in terms of the network lifetime. Below, we discuss our simulation results in detail.

5.2. Fixed triangleRSQ

In this set of experiments, we “x the locations of regions R , S , and query source Q and measure the performance of our algorithms for various values of transmission radii and join-selectivity factors. In particular, we choose coordinates (0, 0), (5, 9.5), and (9.5, 0) for R , Q , and S respectively. The total communication cost incurred by various algorithms for 8000 newly generated tuples of R and S is shown in Fig. 5a...c. We have looked at three transmission radii viz. 0.13, 0.15, and 0.18 units. Lower transmission radii left the sensor network disconnected, and the trend observed for these three transmission radii values was sufficient to infer behavior for larger transmission radii. From Fig. 5a...c, we can see that the Suboptimal Heuristic performs very close to the OptBased Algorithm, and significantly outperforms (upto 100%) the Naive and Centroid Approaches for most parameter values. Sometimes the Suboptimal Heuristic even outperforms the OptBased Algorithm by a small margin.⁸ The performance of the Naive approach worsens drastically with the increase in the join-selectivity factor, since the routing cost of the joined tuples from the join region R (or S) to the query source Q becomes more dominant. For sake of clarity, we have not shown the Naive Approach data points for high join-selectivity factors. Also, note that with the increase in transmission radius and/or selectivity factor, the relative benefit of Suboptimal Heuristic over the Centroid Approach reduces. In particular, for extremely large

⁸ Note that this does not contradict the optimality of the Optimal Algorithm, since the OptBased is only based on the Optimal Algorithm for real sensor networks.

