# Purdue University Purdue e-Pubs

Computer Science Technical Reports

Department of Computer Science

2002

# Evaluating Probabilistic Queries over Imprecise Data

Reynold Cheng

Dmitri V. Kalashnikov

Sunil Prabhakar

Purdue University, sunil@cs.purdue.edu

Report Number: 02-026

Cheng, Reynold; Kalashnikov, Dmitri V.; and Prabhakar, Sunil, "Evaluating Probabilistic Queries over Imprecise Data" (2002). Computer Science Technical Reports. Paper 1544. http://docs.lib.purdue.edu/cstech/1544

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries. Please contact epubs@purdue.edu for additional information.

TR #02-026

# Evaluating Probabilistic Queries over Imprecise Data \*

Reynold Cheng Dmitri V. Kalashnikov Sunil Prabhakar Department of Computer Science, Purdue University Email: {ckcheng,dvk,sunil}@cs.purdue.edu

# **ABSTRACT**

Many applications employ sensors for monitoring entities such as temperature and wind speed. A centralized database tracks these entities to enable query processing. Due to continuous changes in these values and limited resources (e.g., network bandwidth and battery power), it is often infeasible to store the exact values at all times. A similar situation exists for moving object environments that track the constantly changing locations of objects. In this environment, it is possible for database queries to produce incorrect or invalid results based upon old data. However, if the degree of error (or uncertainty) between the actual value and the database value is controlled, we can place more confidence in the answers to queries. More generally, query answers can be augmented with probabilistic estimates of the validity of the answers. In this chapter we study probabilistic query evaluation based upon uncertain data. A classification of queries is made based upon the nature of the result set. For each class, we develop algorithms for computing probabilistic answers. We address the important issue of measuring the quality of the answers to these queries, and provide algorithms for efficiently pulling data from relevant sensors or moving objects in order to improve the quality of the executing queries. Extensive experiments are performed to examine the effectiveness of several data update policies.

# 1. INTRODUCTION

In many applications, sensors are used to continuously track or monitor the status of an environment. Readings from the sensors are sent back to the application, and decisions are made based on these readings. For example, temperature sensors installed in a building are used by a central air-conditioning system to decide whether the temperature of any room needs to be adjusted or to detect other problems. Sensors distributed in the environment can be used to

\* Portions of this work were supported by an Intel Ph.D. fellowship, NSF CAREER grant IIS-9985019, NSF grant 0010044-CCR

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGMOD 2003, June 9-12, San Diego, CA. Copyright 2003 ACM 1-58113-634-X/03/06 ...\$5,00. detect if hazardous materials are present and how they are spreading. In a moving object database, objects are constantly monitored and a central database may collect their undated locations.

The framework for many of these applications includes a database or server to which the readings obtained by the sensors or the locations of the moving objects are sent. Users query this database in order to find information of interest. Due to several factors such as limited network bandwidth to the server and limited battery power of the mobile device, it is often infeasible for the database to contain the exact status of an entity being monitored at every moment in time. An inherent property of these applications is that readings from sensors are sent to the central server periodically. In particular, at any given point in time, the recorded sensor reading is likely to be different from the actual value. The correct value of a sensor's reading is known only when an update is received. Under these conditions, the data in the database is only an estimate of the actual state of the environment at most points in time.

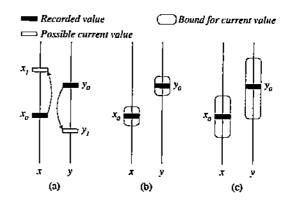


Figure 1: Example of sensor data and uncertainty.

This inherent uncertainty of data affects the accuracy of answers to queries. Figure 1(a) illustrates a query that determines the sensor (either x or y) that reports the lower temperature reading. Based upon the data available in the database ( $x_0$  and  $y_0$ ), the query returns "x" as the result. In reality, the temperature readings could have changed to values  $x_1$  and  $y_1$ , in which case "y" is the correct answer. This example demonstrates that the database does not always truly capture the state of the external world, and the value of the sensor readings can change without being rec-

ognized by the database. Sistla et. al. [14] identify this type of data as a *dynamic attribute*, whose value changes over time even if it is not explicitly updated in the database. In this example, because the exact values of the data items are not known to the database between successive updates, the database incorrectly assumes that the recorded value is the actual value and produces incorrect results.

Given the uncertainty of the data, providing meaningful answers seems to be a futile exercise. However, one can argue that in many applications, the values of objects cannot change drastically in a short period of time; instead, the degree and/or rate of change of an object may be constrained. For example, the temperature recorded by a sensor may not change by more than a degree in 5 minutes. Such information can help solve the problem. Consider the above example again. Suppose we can provide a guarantee that at the time the query is evaluated, the actual values monitored by x and y could be no more than some deviations,  $d_x$  and  $d_y$ , from  $x_0$  and  $y_0$ , respectively, as shown in Figure 1(b). With this information, we can state with confidence that x yields the minimum value.

In general, the uncertainty of the objects may not allow us to identify a single object that has the minimum value. For example, in Figure 1(c), both x and y have the possibility of recording the minimum value since the reading of x may not be lower than that of y. A similar situation exists for other types of queries such as those that request a numerical value (e.g. "What is the lowest temperature reading?"). For these queries too, providing a single value may be infeasible due to the uncertainty in each object's value. Instead of providing a definite answer, the database can associate different levels of confidence with each answer (e.g. as a probability) based upon the uncertainty of the queried objects.

The notion of probabilistic answers to queries over uncertain data has not been well studied. Wolfson et. al briefly touched upon this idea [16] for the case of range queries in the context of a moving object database. The objects are assumed to move in straight lines with a known average speed. The answers to the queries consist of objects' identities and the probability that each object is located in the specified range. In this chapter we extend the notion of probabilistic queries to cover a much broader class of queries. The class of queries considered includes aggregate queries that compute answers over a number of objects. We also discuss the importance of the nature of answer requested by a query (identity of object versus the value). For example, we show that there is a significant difference between the following two queries: "Which object has the minimum temperature?" versus "What is the minimum temperature?". Furthermore, we relax the model of uncertainty so that any reasonable model can be used by the application. Our techniques are applicable to the common models of uncertainty that have been proposed elsewhere.

The probabilities in the answer allow the user to place appropriate confidence in the answer as opposed to having an incorrect answer or no answer at all. Depending upon the application, one may choose to report only the object with the highest probability as having the minimum value, or only those objects whose probability exceeds a minimum probability threshold. Our proposed work will be able to work with any of these models.

Answering aggregate queries (such as minimum or average) is much more challenging than range queries, especially

in the presence of uncertainty. The answer to a probabilistic range query consists of a set of objects along with a non-zero probability that the object lies in the query range. Each object's probability is determined by the uncertainty of the object's value and the query range. However, for aggregate queries, the interplay between multiple objects is critical. The resulting probabilities are greatly influenced by the uncertainty of attribute values of other objects. For example, in Figure 1(c) the probability that x has the minimum value is affected by the relative value and bounds for y.

A probabilistic answer also reflects a certain level of uncertainty that results from the uncertainty of the queries object values. If the uncertainty of all (or some) of the objects was reduced (or eliminated completely), the uncertainty of the result improves. For example, without any knowledge about the value of an object, one could arguably state that it falls within a query range with 50% probability. On the other hand, if the value is known perfectly, one can state with 100% confidence that the object either falls within or outside the query range. Thus the quality of the result is measured by degree of ambiguity in the answer. We therefore need metrics to evaluate the quality of a probabilistic answer. We propose metrics for evaluating the quality of the probabilistic answers in this chapter. As we shall see, it turns out that different metrics are suitable for different classes of queries.

It is possible that the quality of a query result may not be acceptable for certain applications - a more definite result may be desirable. Since the poor quality is directly related to the uncertainty in the object values, one possibility for improving the results is to delay the query until the quality improves. However this is an unreasonable solution due to the increased query response time. Instead, the database could request updates from all objects (e.g. sensors) - this solution incurs a heavy load on the resources. In this chapter, we propose to request updates only from objects that are being queried, and furthermore those that are likely to improve the quality of the query result. We present a number of heuristics and an experimental evaluation. These policies attempt to optimize the use of the constrained resource (e.g. network bandwidth to the server) to improve average query quality.

It should be noted that the imprecision in the query answers is inherent in this problem (due to uncertainty in the actual values of the dynamic attribute), in contrast to the problem of providing approximate answers for improved performance wherein accuracy is traded for efficiency.

To sum up, the contributions introduced in this chapter

- A broad classification of probabilistic queries over uncertain data, based upon a flexible model of uncertainty;
- Techniques for evaluating probabilistic queries, including optimizations;
- Metrics for quantifying the quality of answers to probabilistic queries;
- Policies for improving the quality of answers to probabilistic queries under resource constraints.

The rest of this chapter is organized as follows. In Section 2 we describe a general model of uncertainty, and the

concept of probabilistic queries. Sections 3 and 4 discuss the algorithms for evaluating different kinds of probabilistic queries. Section 5 discusses quality metrics that are appropriate to these queries. Section 6 proposes update policies that improve the query answer quality. We present an experimental evaluation of the effectiveness of these update policies in Section 7. Section 8 discusses related work and Section 9 concludes the chapter.

# 2. PROBABILISTIC QUERIES

In this section, we describe the model of uncertainty considered in this chapter. This is a generic model, as it can accommodate a large number of application peradigms. Based on this model, we introduce a number of probabilistic queries.

# 2.1 Uncertainty Model

One popular model for uncertainty for a dynamic attribute is that at any point in time, the actual value is within a certain bound, d of its last reported value. If the actual value changes further than d, then the sensor reports its new value to the database and possibly changes d. For example, [16] describes a moving-object model where the location of an object is a dynamic attribute, and an object reports its location to the server if its deviation from the last reported location exceeds a threshold. Another model assumes that the attribute value changes with known speed, but the speed may change each time the value is reported. Other models include those that have no uncertainty. For example, in [4], the exact speed and direction of movement of each object are known. This model requires updates at the server whenever an object's speed or direction change.

For the purpose of our discussion, the exact model of uncertainty is unimportant. All that is required is that at the time of query execution the range of possible values of the attribute of interest are known. We are interested in queries over some dynamic attribute, a, of a set of database objects, T. Also, we assume that a is a real-valued attribute, although our models and algorithms can be extended to other domains easily. We denote the ith object of T by  $T_i$  and the value of attribute a of  $T_i$  by  $T_i$ .a (where i = 1, ..., |T|). Throughout this chapter, we treat  $T_{i,a}$  at any time instant t as a continuous random variable. Sometimes instead of writing  $T_{i}.a(t)$  meaning the value of the attribute at time instant t we write just  $T_{i}.a$  for short. The uncertainty of  $T_{i,a}(t)$  can be characterized by the following two definitions (we use pdf to abbreviate the term "probability density function"):

Definition 1: An uncertainty interval of  $T_{i.a}(t)$ , denoted by  $U_i(t)$ , is an interval  $[l_i(t), u_i(t)]$  such that  $l_i(t)$  and  $u_i(t)$  are real-valued functions of t,  $l_i(t) \leq u_i(t)$ , and that the conditions  $u_i(t) \geq l_i(t)$  and  $T_{i.a}(t) \in [l_i(t), u_i(t)]$  always hold.

Definition 2: An uncertainty pdf of  $T_{i.a}(t)$ , denoted by  $f_{i}(x,t)$ , is a pdf of  $T_{i.a}(t)$  such that  $f_{i}(x,t) = 0$  for  $\forall x \notin U_{i}(t)$ .

Notice that since  $f_i(x,t)$  is a pdf, it has the property that  $f_{I_i(t)}^{u_i(t)} f_i(x,t) dx = 1$ . The above definition specifies the uncertainty of  $T_i.a(t)$  in terms of a closed interval and the probability distribution of  $T_i.a(t)$ . Notice that this definition does not specify how the uncertainty interval evolves

over time, and what the pdf  $f_i(x,t)$  is inside the uncertainty interval. The only requirement for  $f_i(x,t)$  is that its value is 0 outside the uncertainty interval. Usually, the scope of uncertainty is determined by the last recorded value, the time elapsed since its last update, and some application-specific assumptions. For example, one may decide that  $U_i(t)$  contains all the values within a distance of  $(t-t_{update}) \times v$  from its last reported value, where  $t_{update}$  is the time that the last update was obtained, and v is the maximum rate of change of the value. One may also specify that  $T_{i.a}(t)$  is uniformly distributed inside the interval, i.e.,  $f_i(x,t) = 1/[u_i(t)-l_i(t)]$  for  $\forall x \in U_i(t)$ , assuming that  $u_i(t) > l_i(t)$ . It should be noted that the uniform distribution represents the worst-case uncertainty over a given interval.

# 2.2 Classification of Probabilistic Oueries

We now present a classification of probabilistic queries and examples of common representative queries for each class. We identify two important dimensions for classifying database queries. First, queries can be classified according to the nature of the answers. An entity-based query returns a set of objects (e.g., sensors) that satisfy the condition of the query. A value-based query returns a single value, examples of which include querying the value of a particular sensor, and computing the average value of a subset of sensor readings. The second property for classifying queries is whether aggregation is involved. We use the term aggregation loosely to refer to queries where interplay between objects determines the result. In the following definitions, we use the following naming convention: the first letter is either E (for entity-based queries) or V (for value-based queries).

# 1. Value-based Non-Aggregate Class

This is the simplest type of query in our discussions. It returns an attribute value of an object as the only answer, and involves no aggregate operators. One example of a probabilistic query for this class is the VSingleQ:

Definition 3: Probabilistic Single Value Query (VSingleQ) When querying  $T_k.a(t)$ , a VSingleQ returns bounds l and u of an uncertainty region of  $T_k.a(t)$  and its pdf  $f_k(x,t)$ .

An example of VSingleQ is "What is the wind speed recorded by sensor  $s_{22}$ ?" Observe how this definition expresses the answer in terms of a bounded probabilistic value, instead of a single value. Also notice that  $\int_{t_i(t)}^{u_i(t)} f_i(x,t) dx = 1$ .

# 2. Entity-based Non-Aggregate Class

This type of query returns a set of objects, each of which satisfies the condition(s) of the query, independent of other objects. A typical example of this class is the ERQ, defined below.

Definition 4: Probabilistic Range Query (ERQ) Given at time instant t a closed interval [l,u], where  $l,u\in\Re$  and  $l\leq u$ , an ERQ returns set R of all tuples  $(T_i,p_i)$ , where  $p_i$  is the non-zero probability that  $T_i.a(t)\in[l,u]$ , i.e.  $R=\{(T_i,p_i):p_i=P(T_i.a(t)\in[l,u])\text{ and }p_i>0\}$ 

An ERQ returns a set of objects, augmented with probabilities, that satisfy the query interval.

## 3. Entity-based Aggregate Class

The third class of query returns a set of objects which satisfy an aggregate condition. We present the definitions of three typical queries for this class. The first two return objects with the minimum or maximum value of  $T_i.a$ :

Definition 5: Probabilistic Minimum (Maximum) Query (EMinQ (EMaxQ)) An EMinQ (EMaxQ) returns set R of all tuples  $(T_i, p_i)$ , where  $p_i$  is the non-zero probability that  $T_i$ , a is the minimum (maximum) value of a among all objects in T.

A one-dimensional nearest neighbor query, which returns object(s) having a minimum absolute difference of  $T_i$ . a and a given value q, is also defined:

Definition 6: Probabilistic Nearest Neighbor Query (ENNQ) Given a value  $q \in \mathbb{R}$ , an ENNQ returns set R of all tuples  $(T_i, p_i)$ , where  $p_i$  is the non-zero probability that  $|T_i \cdot a - q|$  is the minimum among all objects in T.

Note that for all the queries we defined in this class the condition  $\sum_{T_i \in R} p_i = 1$  holds.

## 4. Value-based Aggregate Class

The final class involves aggregate operators that return a single value. Example queries for this class include:

Definition 7: Probabilistic Average (Sum) Query (VAvgQ (VSumQ)) A VAvgQ (VSumQ) returns bounds l and u of an uncertainty interval and pdf  $f_X(x)$  of r.v. X representing the average (sum) of the values of a for all objects in T.

Definition 8: Probabilistic Minimum (Maximum) Value Query (VMinQ (VMaxQ)) A VMinQ (VMaxQ) returns bounds l and u of an uncertainty interval and pdf  $f_X(x)$  of r.v. X representing the minimum (maximum) value of a among all objects in T.

All these aggregate queries return answers in the form of a pdf  $f_X(x)$  and a closed interval [l, u], such that  $\int_l^u f_X(x) dx = 1$ .

Table 2.2 summarizes the basic properties of the probabilistic queries discussed above. For illustrating the difference between probabilistic and non-probabilistic queries, the last row of the table lists the forms of answers expected if probability information is not augmented to the result of the queries e.g., the non-probabilistic version of EMaxQ is a query that returns object(s) with maximum values based only on the recorded values of  $T_i$ .a. It can be seen that the probabilistic queries provide more information on the answers than their non-probabilistic counterparts.

Example. We illustrate the properties of the probabilistic queries with a simple example. In Figure 2, readings of four sensors  $s_1, s_2, s_3$  and  $s_4$ , each with a different uncertainty interval, are being queried at time t. Assume that readings of these sensors  $s_1(t)$ ,  $s_2(t)$ ,  $s_3(t)$ , and  $s_4(t)$  are uniformly distributed on their uncertainty intervals  $[l_1, u_1]$ ,  $[l_2, u_2]$ ,  $[l_3, u_3]$  and  $[l_4, u_4]$ . A VSingleQ applied on  $s_4$  at time t gives us the result:  $l_4$ ,  $u_4$ ,  $f_{s_4}(x) = 1/(u_4 - l_4)$ . When an ERQ (represented by the interval  $\{l, u\}$ ) is invoked at time t to find out how likely each reading is inside [l, u], we see that the reading of  $s_1$  is always inside the interval. It therefore

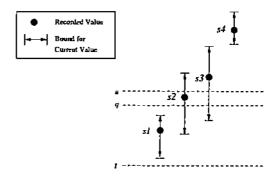


Figure 2: Illustrating the probabilistic queries

has a probability of 1 for satisfying the ERQ. The reading of  $s_4$  is always outside the rectangle, thus it has a probability of 0 of being located inside [l, u]. Since  $U_2(t)$  and  $U_3(t)$  partially overlap [l, u],  $s_2$  and  $s_3$  have some chance of satisfying the query. In this example, the result of the ERQ is:  $\{(s_1, 1), (s_2, 0.7), (s_3, 0.4)\}$ .

In the same figure, an EMinQ is issued at time t. We observe that  $s_1$  has a high probability of having the minimum value, because a large portion of the  $U_1(t)$  has a smaller value than the others. The reading of  $s_1$  has a high chance of being located in this portion because  $s_1(t)$  has the uniform distribution. The reading of  $s_4$  does not have any chance of yielding the minimum value, since none of the values inside  $U_4(t)$  is smaller than others. The result of the EMinQ for this example is:  $\{(s_1, 0.7), (s_2, 0.2), (s_3, 0.1)\}$ . On the other hand, an EMaxQ will return  $\{(s_4, 1)\}$  as the only result, since every value in  $U_4(t)$  is larger than any readings from the other sensors, and we are assured that  $s_4$  yields the maximum value. An ENNQ with a query value q is also shown, where the results are:  $\{(s_1, 0.2), (s_2, 0.5), (s_3, 0.3)\}$ .

When a value-based aggregate query is applied to the scenario in Figure 2, a bounded pdf p(x) is returned. If a VSumQ is issued, the result is a distribution in  $[l_1+l_2+l_3+l_4,u_1+u_2+u_3+u_4]$ ; each x in this interval is the sum of the readings from the four sensors. The result of a VAvgQ is a pdf in  $[(l_1+l_2+l_3+l_4)/4,(u_1+u_2+u_3+u_4)/4]$ . The results of VMinQ and VMaxQ are probability distributions in  $[l_1,u_1]$  and  $[l_4,u_4]$  respectively, since only the values in these ranges have a non-zero probability value of satisfying the queries.

# 3. EVALUATING ENTITY- QUERIES

In this section we examine how the probabilistic entity-based queries introduced in the last section can be answered. We start with the discussion of an ERQ, followed by a more complex algorithm for answering an ENNQ. We also show how the algorithm for answering an ENNQ can be easily changed for EMinQ and EMaxQ.

# 3.1 Evaluation of ERO

Recall that ERQ returns a set of tuples  $(T_i, p_i)$  where  $p_i$  is the non-zero probability that  $T_i.a$  is within a given interval [l, u]. Let R be the set of tuples returned by the ERQ. The algorithm for evaluating the ERQ at time instant t is described in Figure 3.

In this algorithm, each object in T is checked once. To evaluate  $p_i$  for  $T_i$ , we first compute the overlapping interval

Table 1:	Classification	of Probabilistic	Queries.
----------	----------------	------------------	----------

Query Class	Entity-based	Value-based
Aggregate	ENNQ, EMinQ, EMaxQ	VAvgQ, VSumQ, VMinQ, VMaxQ
Non-Aggr.	ERQ	VSingleQ
bilistic)	$\{(T_i, p_i) : 1 \le i \le  T  \land p_i > 0\}$	$l_1u, f_X(x)$
Answer (Non- Prob.)	$\{T_i: 1 \le i \le  T \}$	x ∈ 3ñ

1. 
$$R \leftarrow \emptyset$$
  
2. for  $i \leftarrow 1$  to  $|T|$  do  
(a)  $D \leftarrow U_i(t) \cap [l, u]$   
(b) if  $(D \neq \emptyset)$  then  
i.  $p_i \leftarrow \int_D f_i(x, t) dx$   
ii. if  $p_i \neq 0$  then  $R \leftarrow R \cup \{(T_i, p_i)\}$   
3. return  $R$ 

Figure 3: ERQ Algorithm.

D of the two intervals:  $U_i(t)$  and [l,u] (Step 2a). If D is the empty set, we are assured that  $T_i.a$  does not lie in [l,u], and by the definition of ERQ,  $T_i$  is not included in the result. Otherwise, we calculate the probability that  $T_i.a$  is inside [l,u] by integrating  $f_i(x,t)$  over D, and put the result into R if  $p_i \neq 0$  (Step 2b). The set of tuples  $(T_i,p_i)$ , stored in R, are returned in Step 3.

# 3.2 Evaluation of ENNO

Processing an ENNQ involves evaluating the probability of the attribute a of each object Ti being the closest to (nearest-neighbor of) a value q. In general, this query can be applied to multiple attributes, such as coordinates. In particular, it could be a nearest-neighbor query for moving objects. Unlike the case of ERQ, we can no longer determine the probability for a object independent of the other objects. Recall that an ENNQ returns a set of tuples  $(T_i, p_i)$  where  $p_i$ denotes the non-zero probability that  $T_i$  has the minimum value of  $[T_i, a-q]$ . Let S be the set of objects to be considered by the ENNQ, and let R be the set of tuples returned by the query. The algorithm presented here consists of 4 phases: projection, pruning, bounding and evaluation. The first three phases filter out objects in T whose values of a have no chance of being the closest to q. The final phase, evaluation, is the core of our solution: for every object  $T_i$ that remains after the first three phases, the probability that  $T_i.a$  is nearest to q is computed.

# I. Projection Phase.

In this phase, the uncertainty interval of each  $T_i$  a is computed based on the uncertainty model used by the application. Figure 4(a) shows the last recorded values of  $T_i$  a in S at time  $t_0$ , and the uncertainty intervals are shown in Figure 4(b).

# 2. Pruning Phase.

Consider two uncertainty intervals  $U_1(t)$  and  $U_2(t)$ . If the smallest distance between  $U_1(t)$  and q is larger than the largest distance between  $U_2(t)$  and q, we can immediately conclude that  $T_1$  is not an answer to the ENNQ: even if

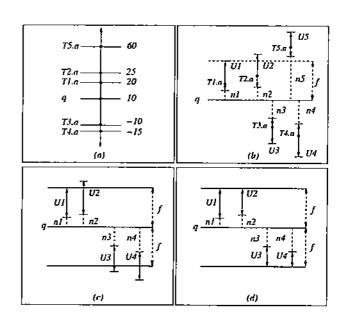


Figure 4: Phases of the ENNQ algorithm.

the actual value of  $T_2.a$  is as far as possible from q,  $T_1.a$  still has no chance to be closer to q than  $T_2.a$ . Based on

```
    for i ← 1 to |S| do

            (a) if q ∈ U<sub>i</sub>(t) then N<sub>i</sub> ← q
            (b) else
            i. if |q - l<sub>i</sub>(t)| < |q - u<sub>i</sub>(t)| then N<sub>i</sub> ← l<sub>i</sub>(t)
            ii. else N<sub>i</sub> ← u<sub>i</sub>(t)
            (c) if |q - l<sub>i</sub>(t)| < |q - u<sub>i</sub>(t)| then F<sub>i</sub> ← u<sub>i</sub>(t)
            (d) else F<sub>i</sub> ← l<sub>i</sub>(t)

    f ← min<sub>1≤i≤|S|</sub>|F<sub>i</sub> - q|; m ← |S|
    for i ← 1 to m do

            if (|N<sub>i</sub> - q| > f) then S ← S \ {T<sub>i</sub>}

    return S
```

Figure 5: Algorithm for the Pruning Phase of ENNQ.

this observation, we can eliminate objects from T by the algorithm shown in Figure 5. In this algorithm,  $N_i$  and  $F_i$  record the closest and farthest possible values of  $T_i.a$  to q, respectively. Steps 1(a) to (d) assign proper values to  $N_i$  and  $F_i$ . If q is inside the interval  $U_i(t)$ , then  $N_i$  is taken as the point q itself. Otherwise,  $N_i$  is either  $l_i(t)$  or  $u_i(t)$ , depending on which value is closer to q.  $F_i$  is assigned in

a similar manner. After this phase, S contains the (possibly fewer) objects which must be considered by q. This is the minimal set of objects which must be considered by the query since any of them can have a value of  $T_i.a$  closest to q. Figure 4(b) illustrates how this phase removes  $T_5$ , which is irrelevant to the ENNQ, from S.

# 3. Bounding Phase.

For each object in S, there is no need to examine all portions in its uncertainty interval. We only need to look at the regions that are located no farther than f from q. We do this conceptually by drawing a bounding interval B of length 2f, centered at q. Any portion of the uncertainty interval outside B can be ignored. Figure 4(c) shows a bounding interval with length 2f, and (d) illustrate the result of this phase.

The phases we have just described attempt to reduce the number of objects to be evaluated, and derive an upper bound on the range of values to be considered.

#### 4. Evaluation Phase.

Based on S and the bounding interval  $B_i$ , our aim is to calculate, for each object in S, the probability that it is the nearest neighbor of q. In the pruning phase, we have already found  $N_i$ , the point in  $U_i(t)$  nearest to q. Let us call  $|N_i - q|$  the near-distance of  $T_i$ , or  $n_i$ . Let as define new r.v.  $X_i$  such

```
    R ← ∅
    Sort the elements in S in ascending order of n<sub>i</sub>, and rename the sorted elements in S as T<sub>1</sub>, T<sub>2</sub>,..., T<sub>|S|</sub>
    n<sub>|S|+1</sub> ← f
    for i ← 1 to |S| do

            (a) p<sub>i</sub> ← 0
            (b) for j ← i to |S| do
            i. p ← ∫<sub>n<sub>j</sub></sub><sup>n<sub>j+1</sub></sup> f<sub>i</sub>(x) · ∏<sub>k=1∧k≠i</sub><sup>j</sup>(1 - F<sub>k</sub>(x)) dx
            ii. p<sub>i</sub> ← p<sub>i</sub> + p
            (c) R ← R ∪ {(T<sub>i</sub>, p<sub>i</sub>)}

    return R
```

Figure 6: Algorithm for the Evaluation Phase of ENNQ.

that  $X_i = |T_i.a(t) - q|$ . Also, let  $F_i(x)$  be the  $X_i$ 's cdf, i.e.,  $F_i(x) = P(|T_i.a(t) - q| \le x)$ , and  $f_i(x)$  be its pdf. Figure 6 presents the algorithm for this phase.

Note that if  $T_{i.a}(t)$  has no uncertainty i.e.,  $U_{i}(t)$  is exactly equal to  $T_{i.a}(t)$ , the evaluation phase algorithm needs to be modified. Our technical report [2] discusses how this algorithm can be changed to adapt to such situations. In the rest of this section, we will explain how the evaluation phase works, assuming non-zero uncertainty.

Evaluation of  $F_i(x)$  and  $f_i(x)$  To understand how the evaluation phase works, it is crucial to know how to obtain  $F_i(x)$ . As introduced before,  $F_i(x)$  is the cdf of  $X_i$ , and thus  $F_i(x) \stackrel{\text{def}}{=} P(|T_i.a(t) - q| \le x)$ . We illustrate the evaluation of  $F_i(x)$  in Figure 7.

Recall that  $f_i(x)$  is the pdf of  $X_i$  and  $f_i(x,t)$  is the pdf of  $T_i.a(t)$ . If  $F_i(x)$  is a differentiable,  $f_i(x)$  is the derivative of  $F_i(x)$ .

Evaluation of  $p_i$ . We can now explain how  $p_i$ , the probability that  $T_{i.a}$  is closest to  $q_i$  is computed. In terms of  $X_i$ 's

- 1. if  $x < n_i$  return 0
- 2. if  $x > |q F_i|$ , return 1
- 3.  $D \leftarrow U_i(t) \cap [q-x,q+x]$
- 4. return  $\int_D f_i(x,t) dx$

Figure 7: Computation of  $F_i(x)$ .

the question is formulated as how  $p_i$ , the probability that  $X_i$  has the minimum value among all  $X_i$ 's, is computed.

Let  $f_i(x) dx$  for indefinitely small dx be the probability that (1)  $X_i \in [x, x+dx]$  and (2)  $X_i = \min_{1 \le k \le |S|} X_k$ . Then Equation 1 outlines the structure of our solution:

$$p_i = \int_{n_i}^f \hat{f}_i(x) dx \tag{1}$$

The  $\hat{f}_i(x) dx$  is equal to probability  $\mathrm{P}(X_i \in [x, x+dx])$  times the probability that  $X_i$  is the minimum among the all  $X_i's$ . The former probability is equal to  $f_i(x) dx$  since  $f_i(x)$  is  $X_i$ 's pdf. The latter probability is equal to the probability that each  $X_j$  in S except for  $X_i$  have values greater than  $X_i$ , which equals to  $\prod_{k=1 \land k \neq i}^{|S|} \mathrm{P}(X_k > x)$ , which also can be written as  $\prod_{k=1 \land k \neq i}^{|S|} \left(1 - F_k(x)\right)$ . Thus the formula for  $p_i$  can be written as:

$$p_{i} = \int_{n_{i}}^{f} f_{i}(x) \cdot \prod_{k=1 \land k \neq i}^{|S|} (1 - F_{k}(x)) dx$$
 (2)

Observe that each  $1 - F_k(x)$  term registers the probability that  $T_k \cdot a$  is forther from q than  $T_i \cdot a$ .

Efficient Computation of  $p_i$  The computation time for  $p_i$  can be improved. Note that  $F_k(x)$  has a value of 0 if  $x < n_k$ . This means if  $x < n_k$  then  $1-F_k(x)$  is always I, and  $T_k$  has no effect on the computation of  $p_i$ . Instead of always considering |S| - 1 objects in the  $\prod$  term of Equation 2 throughout  $[n_i, f]$ , we may actually consider fewer objects in some ranges, resulting in a better computation speed.

This can be achieved by first sorting the objects according to their near\_distance from q. Next, the integration interval  $[n_i, f]$  is broken down into a number of intervals, with end points defined by the near\_distance of the objects. The probability of an object having a value of a closest to q is then evaluated for each interval in a way similar to Equation 2, except that we only consider  $T_k.a$  with non-zero  $F_k(x)$ . Then  $p_i$  is equal to the sum of the probability values for all these intervals. The final formula for  $p_i$  becomes:

$$p_{i} = \sum_{j=i}^{|S|} \int_{n_{j}}^{n_{j+1}} f_{i}(x) \cdot \prod_{k=1, k \neq i}^{j} \left(1 - F_{k}(x)\right) dx \tag{3}$$

Here we let  $n_{|S|+1}$  be f for notational convenience. Instead of considering |S|-1 objects in the  $\prod$  term, Equation 3 only handles j-1 objects in interval  $[n_j, n_j+1]$ . This optimization is shown in Figure 6.

Example Let us use our previous example to illustrate how the evaluation phase works. After 4 objects  $T_1, \ldots, T_4$  were captured (Figure 4(d)), Figure 8 shows the result after these objects have been sorted in ascending order of their near\_distance, with the x-axis being the absolute difference

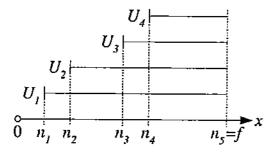


Figure 8: Illustrating the evaluation phase.

of  $T_{i.a}$  from q, and  $n_5$  equals f. The probability  $p_i$  of each  $T_{i.a}$  being the nearest neighbor of q is equal to the integral of  $\hat{f}_i(x)$  over the interval  $[n_i, n_5]$ .

Let us see how we evaluate uncertainty intervals when computing  $p_2$ . Equation 3 tells us that  $p_2$  is evaluated by integrating over  $[n_2, n_5]$ . Since objects are sorted according to  $n_i$ , we do not need to consider all 5 of them throughout  $[n_2, n_5]$ . Instead, we split  $[n_2, n_5]$  into 3 sub-intervals, namely  $[n_2, n_3]$ ,  $[n_3, n_4]$  and  $[n_4, n_5]$ , and consider possibly fewer uncertainty intervals in each sub-interval. For example, in  $[n_2, n_3]$ , only  $U_1$  and  $U_2$  need to be considered.

# 3.3 Evaluation of EMinQ and EMaxQ

We can treat EMinQ and EMaxQ as special cases of ENNQ. In fact, answering an EMinQ is equivalent to answering an ENNQ with q equals the minimum lower bound of all  $U_i(T)$  in T. We can therefore modify the ENNQ algorithm to solve an EMinQ as follows: after the projection phase, we evaluate the minimum value of  $l_i(t)$  among all uncertainty intervals. Then we set q to that value. We then obtain the results to the EMinQ after we execute the rest of the ENNQ algorithm. Solving an EMaxQ is symmetric to solving an EMinQ in which we set q to the maximum of  $u_i(t)$  after the projection phase of ENNQ.

# 4. EVALUATING VALUE- QUERIES

In this section, we discuss how to answer the probabilistic value-based queries defined in Section 2.2.

# 4.1 Evaluation of VSingleO

Evaluating a VSingleQ is simple, since by the definition of VSingleQ, only one object,  $T_k$ , needs to be considered. Suppose VSingleQ is executed at time t. Then the answer returned is the uncertainty information of  $T_k.a$  at time t, i.e.,  $l_k(t)$ ,  $u_k(t)$  and its pdf  $f_k(x,t)$ .

# 4.2 Evaluation of VSumQ and VAvgQ

Let us first consider the case where we want to find the sum of two uncertainty intervals  $\{l_1(t), u_1(t)\}$  and  $\{l_2(t), u_2(t)\}$  for objects  $T_1$  and  $T_2$ . Notice that the values in the answer that have non-zero probability values lie in the range  $[l_1(t) + l_2(t), u_1(t) + u_2(t)]$ . For any x inside this interval,  $f_X(x)$  (the pdf of random variable  $X = T_1.a + T_2.a$ ) is:

$$f_X(x) = \int_{\max\{t_1(t), x - u_2(t)\}}^{\min\{u_1(t), x - t_2(t)\}} f_1(y, t) f_2(x - y, t) \, dy \qquad (4)$$

The lower (upper) bound of the integration interval are evaluated according to the possible minimum (maximum) value of  $T_1.a.$ 

We can generalize this result for summing the uncertainty intervals of |T| objects by picking two intervals, summing them up using the above formula, and using the resulting interval to add to another interval. The process is repeated until we finish adding all the intervals. The resulting interval should have the following form:

$$[\sum_{i=1}^{|T|} l_i(t), \sum_{i=1}^{|T|} u_i(t)]$$

VAvgQ is essentially the same as VSumQ except for a division by the number of objects over which the aggregation is applied.

# 4.3 Evaluation of VMinQ and VMaxO

To answer a VMinQ, we need to find the bounds of uncertainty region [l,u], and pdf  $f_X(x)$  of r.v.  $X=\min_{1\leq i\leq |T|}T_i.a(t)$ . Like for EMinQ the lower bound l can be set as  $\min_{1\leq i\leq |S|}l_i(t)$  and upper bound u as  $\min_{1\leq i\leq |S|}u_i(t)$ , because X cannot take values outside [l,u] The steps of the algorithm are similar to first three phases of ENNQ (projection, pruning, bounding) when q is set to be equal to l. Each  $T_i$  such that  $l_i(t)>u$  is removed from set S of the relevant tuples. Then all tuples in S are sorted in ascending order of  $l_i$ . For notational convenience we introduce an additional parameter  $l_{|S|+1}$  and set it equal to u.

To compute  $f_X(x)$  notice that probability  $P(X \in [x, x + dx])$  for some small dx can be computed as sum of probabilities for each  $T_{i\cdot a}(t)$  to be inside [x, x + dx] times the probability that the other  $T_{i\cdot a}(t)$ 's are in  $(x+dx, +\infty)$ . As we tend dx to zero we have:

$$f_X(x) = \sum_{i=1}^{|S|} \left( f_i(x,t) \cdot \prod_{k=1}^{|S|} \left( 1 - F_k(x,t) \right) \right), \ \forall x \in [l,u] \ (5)$$

Since if  $x \in [l_j, l_{j+1}]$  and k > j, then  $F_k(x, t) = 0$ , therefore terms  $(1 - F_k(x, t))$  are equal to 1 for such x's and k's and need not be considered by the formula. The simplified formula is thus:

$$f_X(x) = \sum_{i=1}^{|S|} \left( f_i(x, t) \cdot \prod_{k=1 \land k \neq 1}^{j} (1 - F_k(x, t)) \right), \ \forall x \in [l_j, l_{j+1}], 1 \le j \le |S|$$

Also notice that if  $x \in [l_j, l_{j+1}]$  and k > j, then  $f_i(x, t) = 0$ . Thus formula for  $f_X(x)$  can be written as:

$$f_X(x) = \sum_{i=1}^{j} \left( f_i(x,t) \cdot \prod_{k=1 \land k \neq i}^{j} (1 - F_k(x,t)) \right), \ \forall x \in [l_j, l_{j+1}], 1 \le j \le |S|$$
(7)

VMaxQ is handled in an analogous fashion.

# 5. QUALITY OF PROBABILISTIC RESULTS

In this section, we discuss several metrics for measuring the quality of the results returned by probabilistic queries. It is interesting to see that different metrics are suitable for different query classes.

# 5.1 Entity-Based Non-Aggregate Queries

For queries that belong to the entity-based non-aggregate query class, it suffices to define the quality metric for each  $(T_i, p_i)$  individually, independent of other tuples in the result. This is because whether an object satisfies the query or not is independent of the presence of other objects. We illustrate this point by explaining how the metric of ERQ is defined.

For an ERQ with query range [l,u], the result is the best if we are sure either  $T_{i\cdot a}$  is completely inside or outside [l,u]. Uncertainty arises when we are less than 100% sure whether the value of  $T_{i\cdot a}$  is inside [l,u]. We are confident that  $T_{i\cdot a}$  is inside [l,u] if a large part of  $U_i(t)$  overlaps [l,u] i.e.,  $p_i$  is large. Likewise, we are also confident that  $T_{i\cdot a}$  is outside [l,u] if only a very small portion of  $U_i(t)$  overlaps [l,u] i.e.,  $p_i$  is small. The worst case happens when  $p_i$  is 0.5, where we cannot tell if  $T_{i\cdot a}$  satisfies the range query or not. Hence a reasonable metric for the quality of  $p_i$  is:

$$\frac{|p_i - 0.5|}{0.5} \tag{8}$$

In Equation 8, we measure the difference between  $p_i$  and 0.5. Its highest value, which equals 1, is obtained when  $p_i$  equals 0 or 1, and its lowest value, which equals 0, occurs when  $p_i$  equals 0.5. Hence the value of Equation 8 varies between 0 to 1, and a large value represents good quality. Let us now define the *score* of an ERQ:

Score of an ERQ = 
$$\frac{1}{|R|} \sum_{i \in R} \frac{|p_i - 0.5|}{0.5}$$
 (9)

where R is the set of tuples  $(T_i.a, p_i)$  returned by an ERQ. Essentially, Equation 9 evaluates the average over all tuples in R.

Notice that in defining the metric of ERQ, Equation 8 is defined for each  $T_i$ , disregarding other objects. In general, to define quality metrics for the entity-based non-aggregate query class, we can define the quality of each object individually. The overall score can then be obtained by averaging the quality value for each object.

# 5.2 Entity-Based Aggregate Queries

Contrary to an entity-based non-aggregate query, we observe that for an entity-based aggregate query, whether an object appears in the result depends on the existence of other objects. For example, consider the following two sets of answers to an EMinQ:  $\{(T_1.a, 0.6), (T_2.a, 0.4)\}$  and  $\{(T_1.a, 0.6), (T_2.a, 0.3), (T_3.a, 0.1)\}$ . How can we tell which answer is better? We identify two important components of quality for this class: entropy and interval width.

Entropy. Let r.v. X take values from set  $\{x_1, \ldots, x_n\}$  with respective probabilities  $p(x_1), \ldots, p(x_n)$  such that  $\sum_{i=1}^n p(x_i) = 1$ . The entropy of X is a measure H(X):

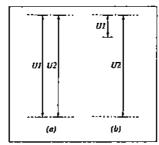
$$H(X) = \sum_{i=1}^{n} p(x_i) \log_2 \frac{1}{p(x_i)}$$
 (10)

If  $x_i$ 's are treated as messages and  $p(x_i)$ 's as their probability to appear, then the entropy H(X) measures the average number of bits required to encode X, or the amount of information carried in X [13]. If H(X) equals 0, there exists some i such that  $p(x_i) = 1$ , and we are certain that  $x_i$  is the message, and there is no uncertainty associated with X. On the other hand, H(X) attains the maximum value when all the messages are equally likely, in which case H(X) equals  $\log_2 n$ .

Recall that the result to the queries we defined in this class is returned in a set R consisting of tuples  $(T_i, p_i)$ . Let r.v.

Y take value i with probability  $p_i$  if and only if  $(T_i, p_i) \in R$ . The property that  $\sum_{i=1}^{n} p_i = 1$  holds. Then H(Y) measures the uncertainty of the answer to these queries; the lower the value of H(Y), the more certain is the answer.

Bounding Interval. Uncertainty of an answer also depends on another important factor: the bounding interval B. Recall that before evaluating one of these aggregate queries, we need to find B that dictates all possible values we have to consider. Then we consider all the portions of uncertainty intervals that lie within B. Note that the decision of which object satisfies the query is only made within this interval. Also notice that the width of B is determined by the width of the uncertainty intervals associated with objects; a large width of B is the result of large uncertainty intervals. Therefore, if B is small, it indicates that the uncertainty intervals of objects that participate in the final result of the query are also small. In the extreme case, when the uncertainty intervals of participant objects have zero width, then the width of B is zero too. The width of B therefore gives us a good indicator of how uncertain a query answer is.



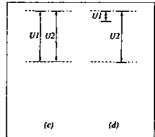


Figure 9: Illustrating how the entropy and the width of B affect the quality of answers for entity-based aggregate queries. The four figures show the uncertainty intervals  $(U_1(t_0))$  and  $U_2(t_0)$  inside B after the bounding phase. Within the same bounding interval, (b) has a lower entropy than (a), and (d) has a lower entropy than (c). However, both (c) and (d) have less uncertainty than (a) and (b) because of smaller bounding intervals.

An example at this point will make our discussions clear. Figure 9 shows four different scenarios of two uncertainty intervals,  $U_1(t_0)$  and  $U_2(t_0)$ , after the bounding phase for an EMinQ. We can see that in (a),  $U_1(t_0)$  is the same as  $U_2(t_0)$ . If we assume a uniform distribution for both uncertainty intervals, both  $T_1$  and  $T_2$  will have equal probability of having the minimum value of a. In (b), it is obvious that  $T_2$  has a much greater chance than  $T_1$  to have the minimum value of a. Using Equation 10, we can observe that the answer in (b) enjoys a lower degree of uncertainty than (a). In (c) and (d), all the uncertainty intervals are halved of those in (a) and (b) respectively. Hence (d) still has a lower entropy value than (c). However, since the uncertainty intervals in (c) and (d) are reduced, their answers should be more certain than those of (a) and (b). Notice that the widths of B for (c) and (d) are all less than (a) and (b).

The quality of entity-based aggregate queries is thus decided by two factors: (1) entropy H(Y) of the result set, and (2) width of B. Their scores are defined as follows:

Score of an Entity, Aggr Query = -H(Y) width of B

Notice that the query answer gets a high score if either H(Y) is low, or the width of B is low. In particular, if either H(Y) or the width of B is zero, then H(Y) = 0 is the maximum score.

# 5.3 Value-Based Oueries

Recall that the results returned by value-based queries are all in the form an uncertainty interval [l, u], and pdf f(x). To measure the quality of such queries, we can use the concept of entropy of a continuous distribution, defined as follows:

$$\hat{H}(X) = -\int_{t}^{u} f(x) \log_{2} f(x) dx$$
 (12)

where  $\hat{H}(X)$  is the entropy of continuous random variable X with pdf f(x) defined in the interval [l,u] [13]. Similar to the notion of entropy,  $\hat{H}(X)$  measures the uncertainty associated with the value of X. Moreover, X attains the maximum value,  $log_2(u-l)$  when X is uniformly distributed in [l,u]. Entropy  $\hat{H}(X)$  can be negative, e.g. for uniform r.v.  $X \sim U[0,\frac{1}{2}]$ .

We use the notion of entropy of a continuous distribution to measure the quality of value-based queries. Specifically, we apply Equation 12 to f(x) as a measure of how much uncertainty is inherent to the answer of a value-based query. The lower the entropy value, the more certain is the answer, and hence the better quality is the answer. We now define the score of a probabilistic value-based query:

Score of a Value-Based Query = 
$$-\hat{H}(X)$$
 (13)

The quality of a value-based query can thus be measured by the uncertainty associated with its result: the lower the uncertainty, the higher score can be obtained as indicated by Equation 13.

Please notice that though not presented here many more different metrics from those discussed in these research are possible; e.g. one might choose the standard deviation as a metric for the quality of VMinQ etc.

# 6. IMPROVING ANSWER QUALITY

In this section, we discuss several update policies that can be used to improve the quality of probabilistic queries, defined in the last section. We assume that the sensors coperate with the central server i.e., a sensor can respond to update requests from the sensor by sending the newest value to the server, as in the system model described in [10].

Suppose after the execution of a probabilistic query, some slack time is available for the query. The server can improve the quality of the answers to that query by requesting updates from sensors, so that the uncertainty intervals of some sensor data are reduced, potentially resulting in an improvement of the answer quality. Ideally, a system can demand updates from all sensors involved in the query; however, this is not practical in a limited-bandwidth environment. The issue is, therefore, to improve the quality with as few updates as possible. Depending on the types of queries, we propose a number of update policies.

Improving the Quality of ERQ The policy for choosing objects to update for an ERQ is very simple: choose the

object with the minimum value computed in Formula 8, with an attempt to improve the score of ERQ.

Improving the Quality of Other Queries Several update policies are proposed for queries other than ERO:

- Glb\_RR. This policy updates the database in a roundrobin fashion using the available bandwidth i.e., it updates the data items one by one, making sure that each item gets a fair chance of being refreshed.
- Loc\_RR. This policy is similar to Glb\_RR, except that the round-robin policy is applied only to the data items that are related to the query, e.g., the set of objects with uncertainty intervals overlapping the bounding interval of an EMinQ.
- MinMin. An object with its lower bound of the uncertainty interval equal to the lower bound of B is chosen for update. This attempts to reduce the width of B and improve the score.
- MaxUnc. This heuristic simply chooses the uncertainty interval with the maximum width to update, with an attempt to reduce the overlapping of the uncertainty intervals.
- 5. MinExpEntropy. Another heuristic is to check, for each T<sub>i</sub>.a that overlaps B, the effect to the entropy if we choose to update the value of T<sub>i</sub>.a. Suppose once T<sub>i</sub>.a is updated, its uncertainty interval will shrink to a single value. The new uncertainty is then a point in the uncertainty interval before the update. For each value in the uncertainty interval before the update, we evaluate the entropy, assuming that U<sub>i</sub>(t) shrinks to that value after the update. The mean of these entropy values is then computed. The object that yields the minimum expected entropy is updated.

# 7. EXPERIMENTAL RESULTS

In this section, we experimentally study the relative behaviors of the various update policies described above, with respect to improving the quality of the query results. We will discuss the simulation model followed by the results.

# 7.1 Simulation Model

The evaluation is conducted using a discrete event simulation representing a server with a fixed network bandwidth (B messages per second) and 1000 sensors. Each update from a sensor updates the value and the uncertainty interval for the sensor stored at the server. The uncertainty model used in the experiments is as follows: An update from sensor  $T_i$  at time  $t_{update}$  specifies the current value of the sensor,  $T_{i.a_{srv}}$ , and the rate,  $T_{i.r_{srv}}$  at which the uncertainty region (centered at  $T_{i.a_{srv}}$ ) grows. Thus at any time instant,  $t_i$  following the update, the uncertainty interval  $\{U_i(t)\}$  of sensor  $T_i$  is given by  $T_{i.a_{srv}} \pm T_{i.r_{srv}} \times (t - T_{i.t_{update}})$ . The distribution of values within this interval is assumed to be uniform.

The actual values of the sensors are modeled as random walks within the normalized domain as in [10]. The maximum rate of change of individual sensors are uniformly distributed between 0 and  $R_{max}$ . At any time instant, the value of a sensor lies within its current uncertainty interval specified by the last update sent to the server. An update

from the sensor is necessitated when a sensor is close to the edge of its current uncertainty region. Additionally, in order to avoid excessively large levels of uncertainty, an update is sent if either the total size of the uncertainty region or the time since the last update exceed threshold values.

The representative experiments presented considered either EMinQ or VMinQ queries only. In each experiment the queries arrive at the server following a Poisson distribution with arrival rate  $\lambda_q$ . Each query is executed over a subset of the sensors. The subsets are selected randomly following the 80-20 hot-cold distribution (20% of the sensors are selected 80% of the time). The cardinality of each set was fixed at  $N_{aub}=100$ . The maximum number of concurrent queries was limited to  $N_q=10$ . Each query is allowed to request at most  $N_{msg}$  updates from sensors in order to improve the quality of its result.

In order to study different aspects of the policies, query termination can be specified either as (i) a fixed time interval ( $T_{active}$ ) after which the query is completed even if its requested updates have not arrived (due to network congestion) or (ii) when a target quality ( $\mathcal{G}$ ) is achieved. Depending upon the policy, we study either the average achieved quality (score), the average size of the uncertainty region, or the average response time needed to achieve the desired quality. All measurements were made after a suitable warm up period had elapsed. For fairness of comparison, in each experiment, the arrival of queries as well as the changes to the sensor values was identical.

Table 7.1 summarizes the major parameters and their default values. The simulation parameters were chosen such that average cardinality of the result sets achieved by the best update policies was between 3 and 10.

Table 2: Simulation parameters and their default values

Param	Default	Meaning
$\mathcal{D}$	[0, 1]	Domain of attribute a
Rmax	0.1	Maximum rate of change of a (sec-1)
$N_{q}$	10	Maximum # of concurrent queries
$\lambda_q$	20	Query arrival rate (query/sec)
Naub	100	Cardinality of query subset
Tactive	5	Query active time (sec)
В	350	Network bandwidth (msg/sec)
$N_{msg}$	5	Maximum # of updates per query
Ncone	1	The # of concurrent updates per query

# 7.2 Results

Due to limited space, we only show the most important experimental results. Interested readers are referred to our technical report [2] for more detailed discussions of our experiments. All figures in this section show averages.

Bandwidth. Figure 10 shows scores for EMinQ achieved by various update policies for different values of bandwidth. The quality metric in this case is negated entropy times the size of the uncertainty region of the result set. Figure 11 is analogous to Figure 10 but shows scores for VMinQ instead of EMinQ. The score for VMinQ queries is negated continuous entropy.

In Figures 10 and 11, the scores increase as bandwidth

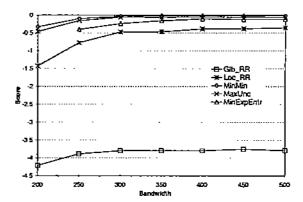


Figure 10: EMinQ score as function of B

increases for all policies, approaching the perfect score of zero for EMinQ. This is explained by the fact that with higher bandwidth the updates requested by the queries are received faster. Thus for higher bandwidth the uncertainty regions for freshly updated sensors tend to be smaller than those using lower bandwidth. Smaller uncertainty regions translate into smaller uncertainty of the result set, and consequently higher score. The reduction in uncertainty regions with increasing bandwidth can be observed from Figure 12.

All schemes that favor updates for sensors being queried significantly outperform the the only scheme that ignores this information: Glb\_RR. The best performance is achieved by the MinMin policy, which updates a sensor with the lower bound of the uncertainty region  $l_i(t)$  equal to the minimum lower bound among all sensors considered by the query. The MinExpEntropy policy showed worse results<sup>1</sup> than the Min-Min and MaxUnc policies in Figures 10 and 12 and worse results than those of the MinMin policy for VMinQ queries, Figure 11. When comparing the MinMin and MaxUnc policies, the better score of the MinMin policy is explained by the fact that the sensor picked for an update by the Min-Min policy tends to have large uncertainty too - in fact, the uncertainty interval is at least as large as the width of the bounding interval. In addition the value of its attribute a tends to have higher probability of being minimum.

Response Time. Figure 13 shows response time as a function of available bandwidth for EMinQ. Unlike the other experiments, in this experiment a query execution is stopped as soon as the goal score  $\mathcal{G}$  (-0.06) is reached. Once again the MinMin strategy showed the best results, reaching the goal score faster than the other policies. The difference in response time is especially noticeable for smaller values of bandwidth, where it is almost twice as good as the other strategies. Predictably, the response time decreases when more bandwidth becomes available.

Arrival Rate. Figures 14 and 15 show the scores achieved by EMinQ and VMinQ queries for various update policies as a function of query arrival rate  $\lambda_q$ . As  $\lambda_q$  increases from 5 to 25, more queries request updates and reduce the uncertainty regions. As a result, the uncertainty decreases, which leads to better scores (Figure 16). When  $\lambda_q$  reaches 25 the entire network bandwidth is utilized. As  $\lambda_q$  continue to increase

<sup>&</sup>lt;sup>1</sup>The experiment with bandwidth of 200 did not complete in time for the submission. The final version of the paper will contain all results.

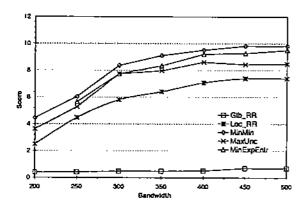


Figure 11: VMinQ score as function of B

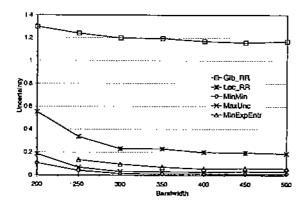


Figure 12: Uncertainty as function of B

queries are able to send fewer requests for updates and receive fewer updates in time, leading to poor result quality and larger uncertainty.

We can observe from Figures 14, 15, and 16 that the relative performance of the various policies remains the same over a wide range of arrival rates  $\{\lambda_q \in [5,45]\}$ .

The experiments show that all policies that favor querybased updates achieve much higher levels of quality. For the queries considered, the MinMin policy gives the best performance. Evaluation of the policies for all types of queries is beyond the scope of this thesis. We plan to address this issue as part of future work.

# 8. RELATED WORK

Many studies have focussed on providing approximate answers to database queries. These techniques approximate query results based only upon a subset of data. In [15], Vrbsky et. al studied how to provide approximate answers to set-valued queries (where a query answer contains a set of objects) and single-valued queries (where a query answer contains a single value). An exact answer E can be approximated by two sets: a certain set C which is the subset of E, and a possible set P such that  $C \cup P$  is a superset of E. Unlike our assumptions, their model assumes there is no uncertainty in the attribute values. Other techniques use precomputation [11], sampling [5] and synopses [1] to produce statistical results. While these efforts investigate ap-

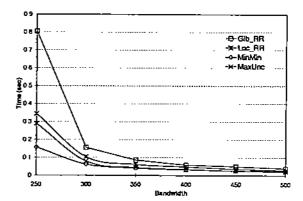


Figure 13: Response time as function of B

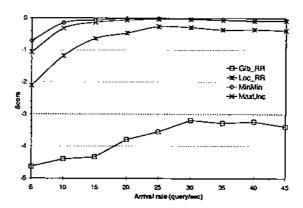


Figure 14: EMinQ score as function of  $\lambda_q$ 

proximate answers based upon a subset of the (exact) values of the data, our work addresses probabilistic answers based upon all the (imprecise) values of the data.

The problem of balancing the tradeoff between precision and performance for querying replicated data was studied by Olston et. al. [9, 8, 10]. In their model, the cache in the server cannot keep track of the exact values of sensor sources due to limited network bandwidth. Instead of storing the actual value for each data item in the server's cache, they propose to store an interval for each item within which the current value must be located. A query is then answered by using these intervals, together with the actual values fetched from the sources. In [9], the problem of minimizing the update cost within an error bound specified by aggregate queries is studied. In [8], algorithms for tuning the intervals of the data items stored in the cache for best performance are proposed. In [10], the problem of minimizing the divergence between the server and the sources given a limited amount of bandwidth is discussed.

Khanna et. al [7] extend Olston's work by proposing an online algorithm that identifies a set of elements with minimum update cost so that a query can be answered within an error bound. Three models of precision are discussed: absolute, relative and rank. In the absolute (relative) precision model, an answer a is called  $\alpha$ -precise if the actual value v deviates from a by not more than an additive (multiplicative) factor of  $\alpha$ . The rank precision model is used to deal

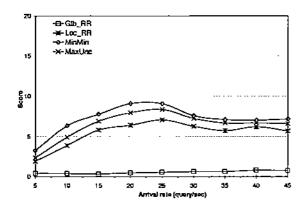


Figure 15: VMinQ score as function of  $\lambda_a$ 

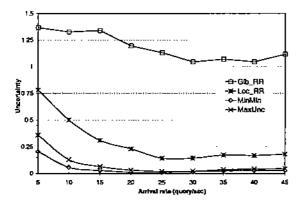


Figure 16: Uncertainty as function of  $\lambda_a$ 

with selection problems which identifies an element of rank r: an answer a is called  $\alpha$ -precise if the rank of a lies in the interval  $[r - \alpha, r + \alpha]$ .

In all the works that we have discussed, the use of probability distribution of values inside the uncertainty interval as a tool for quantifying uncertainty has not been considered. Discussions of queries on uncertainty data were often limited to the scope of aggregate functions. In contrast, our work adopts the notion of probability and provides a paradigm for answering general queries involving uncertainty. We also define the quality of probabilistic query results which, to the best of our knowledge, has not been addressed.

With the exception of [16], we are unaware of any work that discusses the evaluation of a query answer in probabilistic form. The study in [16] is limited to range queries for objects moving in straight lines in the context of a moving-object environment. We extend their ideas significantly by providing probabilistic guarantees to general queries for a generic model of uncertainty. Other related work include [12, 3, 6].

# 9. CONCLUSIONS

In this chapter we studied the problem of augmenting probability information to queries over uncertain data. We propose a flexible model of uncertainty, which is defined by (1) an lower and upper bound, and (2) a pdf of the values inside the bounds. We then explain, from the viewpoint of

a probabilistic query, we can classify queries in two dimensions, based on whether they are aggregate/non-aggregate queries, and whether they are entity-based/value-based. Algorithms for computing typical queries in each query class are demonstrated. We present novel metrics for measuring quality of answers to these queries, and also discuss several update heuristics for improving the quality of results. The benefit of query-based updates was also shown experimentally.

# 10. REFERENCES

- S. Acharya, P. B. Gibbons, V. Poosala, and S. Ramaswamy. Join synopses for approximate query answering. In SIGMOD 1999.
- [2] R. Cheng, D. V. Kalashnikov, and S. Prabhakar. Evaluating probabilistic queries over imprecise data. Technical Report TR 02-026, Department of Computer Science, Purdue University, West Lafayette, IN 47907, Nov. 2002.
- [3] R. Cheng, S. Prabhakar, and D. Kalashnikov. Querying imprecise data in moving object environments. In ICDE'03, Proc. of IEEE Int'l Conf. on Data Engineering, Mar 5-8 2003.
- [4] D.Pfoser and C.S.Jensen. Querying the trajectories of on-line mobile objects. In *MobiDE 2001*, pages 66-73.
- [5] P. B. Gibbons and Y. Matias. New sampling-based summary statistics for improving approximate query answers. In SIGMOD 1998.
- [6] D. Kalashnikov, S. Prabhakar, S. Hambrusch, and W. Aref. Efficient evaluation of continuous range queries on moving objects. In DEXA'02, Proc. of Int'l Conf. on Database and Expert Systems Applications, Sep 2-6 2002.
- [7] S. Khanna and W. Tan. On computing functions with uncertainty. In PODS 2001.
- [8] C. Olston, B. T. Loo, and J. Widom. Adaptive precision setting for cached approximate values. In ACM SIGMOD, 2001.
- [9] C. Olston and J. Widom. Offering a precision-performance tradeoff for aggregation queries over replicated data. In VLDB, 2000.
- [10] C. Olston and J. Widom. Best-effort cache synchronization with source cooperation. In ACM SIGMOD, pages 73-84, 2002.
- [11] V. Poosala and V. Ganti. Fast approximate query answering using precomputed statistics. In *Proc of the* 15th ICDE, page 252, 1999.
- [12] S. Prabhakar, Y. Xia, D. Kalashnikov, W. Aref, and S. Hambrusch. Query indexing and velocity constrained indexing: Scalable techniques for continuous queries on moving objects. *IEEE Transactions on Computers*, 51(10):1124-1140, Oct. 2002.
- [13] C. E. Shannon. The Mathematical Theory of Communication. University of Illinois Press, 1949.
- [14] P. A. Sistla, O. Wolfson, S. Chamberlain, and S. Dao. Querying the uncertain position of moving objects. In Temporal Databases: Research and Practice, number 1399, 1998.
- [15] S. V. Vrbsky and J. W. S. Liu. Producing approximate answers to set- and single-valued queries. The Journal of Systems and Software, 27(3), 1994.

[16] O. Wolfson, P. A. Sistla, S. Chamberlain, and Y. Yesha. Updating and querying databases that track mobile units. *Distributed and Parallel Databases*, 7(3):257-387, 1999.