Global Journals $ensuremath{\mathbb{E}} T_{\ensuremath{\mathbb{E}}} X$ JournalKaleidoscope
TM

Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

Data Aggregators Dr. P. Prabhakar¹ and S. Nageswara Rao² M.Tech Student, Dept of CSE, Madanapalle Institute of Technology and Science, AP, Received: 15 December 2011 Accepted: 31 December 2011 Published: 15 January 2012

7 Abstract

Typically a user desires to obtain the value of some aggregation function over distributed data 8 items. We present a low-cost, scalable technique to answer continuous aggregation queries 9 using a network of aggregators of dynamic data items. In such a network of data aggregators, 10 each data aggregator serves a set of data items at specific coherencies. Our technique involves 11 decomposing a client query into sub-queries and executing sub-queries on judiciously chosen 12 data aggregators with their individual subquery incoherency bounds. We provide a technique 13 for getting the optimal set of sub-queries with their incoherency bounds, which satisfies client 14 query?s coherency requirement with least number of refresh messages sent from aggregators to 15 the client. For estimating the number of refresh messages, we build a query cost model which 16 can be used to estimate the number of messages required to satisfy the client specified 17 incoherency bound. Performance results using real-world traces show that our cost based 18 query planning leads to queries being executed using less than one third the number of 19 messages required by existing schemes. 20

21

Index terms— Content distribution network, continuous query, online decision making, data dissemination,
 coherency, performance.

24 1 INTRODUCTION

25 pplication such as auctions, personal portfolio for financial decisions, sensors based monitoring, route planning based on traffic information, etc., make extensive use of dynamic data. For such applications, data from one or 26 more independent data sources may be aggregated to determine if some action is warranted. Given the increasing 27 number of such applications that make use of highly dynamic data, there is significant interest in systems that 28 can efficiently deliver the relevant updates automatically. Many data intensive applications delivered over the 29 Web suffer from performance and scalability issues. Content distribution networks (CDNs) solved the problem 30 for static content using caches at the edge nodes of the networks. CDNs continue to evolve to serve more 31 and more dynamic applications [1,2]. A dynamically generated web page is usually assembled using a number 32 of static or dynamically generated fragments. The static fragments are served from the local caches whereas 33 dynamic fragments are created either by using the cached data Author? : M.Tech Student, Dept of CSE, 34 35 Madanapalle Institute of Technology and Science, AP, India. E-mail : prabhakar.mth@gmail.com Author ? 36 : Assistant Professor, Dept of CSE, Madanapalle Institute of Technology and Science, AP, India. E-mail : 37 nag_sirisala@yahoo.com or by fetching the data items from the origin data sources. One important question for satisfying client requests through a network of nodes is how to select the best node(s) to satisfy the request. 38 For static pages content requested, proximity to the client and load on the nodes are the parameters generally 39 used to select the appropriate node [3]. In dynamic CDNs, while selecting the nodes node(s) to satisfy the client 40 request, the central site (top-level CDN node) has to ensure that page/data served meets client's coherency 41 requirements also. Techniques to efficiently serve fast changing data items with guaranteed incoherency bounds 42 have been proposed in the literature [4,5]. Such dynamic data dissemination networks can be used to disseminate 43

44 data such as stock quotes, temperature data from sensors, traffic information, and network monitoring data. In

this paper we propose a method to efficiently answer aggregation queries involving such data items. In data dissemination schemes proposed in literature [4,11], a hierarchical network of data aggregators is employed such

that each data aggregator serves the data item at some guaranteed incoherency bound. Incoherency of a data

item at a given node is defined as the difference in value of the data item at the data source and the value at

49 that node. Although CDNs use page-purge [8] based coherency management, we assume that in dynamic data

 $_{\rm 50}$ $\,$ dissemination networks, these messages carry the new data values thereby an invalidation message becomes a

⁵¹ refresh message. For maintaining a certain incoherency bound, a data aggregator gets data updates from the data

52 source or some higher level data aggregator so that the data incoherency is not more than the data incoherency 53 bound. In a hierarchical data dissemination network a higher level aggregator guarantees a tighter incoherency

54 bound compared to a lower level aggregator. Thus, data refreshes are pushed from the data sources to the clients

55 through the network of aggregators.

Data incoherency: data accuracy can be specified in terms of incoherency of a data item, defined as the absolute difference in value of the data item at the data source and the value known to a client of the data. Let ? i (t) denote the value of i th data item at the data source at time t; and let the value the data item known to the client be u i (t). Then the data incoherency at the client is given by |u i (t)-? i (t)|. For a data item which

needs to be refreshed at an incoherency bound C aA (D D D D)

data refresh message is sent to the client as soon as data exceeds C, i.e., |u i (t) -? i (t)| > C.

Network of data aggregators: Data aggregators are one kind of secondary server it serves as data sources (data items). The data refreshes can be done using two mechanisms.(a)Push based mechanism data source send update messages to client on their own.(b)Pull based mechanism data sources send messages to the client only when client makes a request. For scalable handling of push based data dissemination, network of DA's are proposed in the literature [12,15,16].

67 2 Sources

68 Clients Figure 1 : Data dissemination network for multiple data items

In such network of DA's, data refreshes occur from data sources to the client through one or more DA's. In this

paper we assume that each DA maintains its configured incoherency bounds for various data items. Dissemination networks for various data items can be overlaid over a single network of data aggregators as shown in Figure 1.

Thus, From a data dissemination capability point of view, each data aggregator (DA) is characterized by a set

 72 of (d i , c i) pairs, where d i is a data item which the DA can disseminate at an incoherency bound C i .The

configured incoherency bound of data item at a DA can be maintained using any of following methods: (a) the

⁷⁵ data source refreshes the data value of the DA whenever DA's incoherency bound is about to get violated. This

⁷⁶ method is scalability problems. (b) data aggregators with tighter incoherency bound help the DA to maintain

its incoherency bound in scalable manner as explained in [4,7].

Example 1: In a network of data aggregators managing data items D1-D4, various aggregators can be represented as-A1: $\{(D1, 0.5), (D3, 0.2)\}$ A2: $\{(D1, 1.0), (D2, 0.1), (D4, 0.2)\}$

Aggregator A1 can serve values of D1 with an incoherency bound greater than or equal to 0.5 whereas A2 can disseminate the same data item at a looser incoherency bound of 1.0 or more. Usually, client is interested in an aggregation of these dynamic data items at a certain incoherency bound. These continuous queries are used to monitor changes in dynamic data and provide results useful for online decision-making.

For generating the result of a query, data from multiple sources is required. As a result, the query has to 84 be evaluated either at data aggregators or at the client. In this work, our aim is to satisfy the client's query 85 requirements while minimizing the query execution cost in terms of number of dissemination messages. Towards 86 that end, we have achieved the following: 1. Developed techniques for estimating the cost of disseminating a 87 data item, at specified incoherency bound. 2. Using the estimated data dissemination cost, we developed query 88 cost model for estimating the cost of executing an incoherency bounded continuous query. 3. Used the query 89 cost model for assigning a client query to one or more data aggregators so that the query can be executed with 90 the least number of messages. 91

Our work involves dividing the client query into sub-queries and allocating it to different data aggregators for optimal execution. In comparison, all the related work in literature [3,5] propose getting individual data items from the aggregators which, as we show in this report, leads to large number of dissemination messages. In the rest of the Introduction, we present basic concepts underlying incoherency bounded continuous query execution

96 using a distributed network of data aggregators.

⁹⁷ 3 a) Problem Statement and Contributions

In this paper, we develop query cost model for aggregation query involving multiple data items: -Additive aggregation with each data item possibly different weights, and -MIN/MAX aggregation queries. The weighted aggregation query can be mathematically written as: () () 1 q q q i i i n s i v t s t w = = \times ?(1)

V s q is the value of a client query q involving n q data items with the weight of the i th data item being w q i, 1 < i < n q. s i (t) is the value of the i th data item at the data source at time t. Such a query encompasses SQL aggregation operators SUM and AVG besides general weighted aggregation queries such as portfolio queries, involving aggregation of stock prices, weighted with number of shares of stocks in the portfolio. Due to ($\rm D~D~D$ $_{105}$ $\rm~D$)

Year space limitations, we are not presenting execution schemes for other aggregation queries such as MIN/MAX. Interested readers are referred to [13] for the extended version of this paper.

Let the value of i th data item, in Equation (1), known to the client/DA be d i (t). Then the data incoherency is given by |s i (t)-d i (t)|. For a data item which needs to be disseminated at an incoherency bound C the data refresh is sent to the client or lower level DA, if the |s i (t)d i (t)| is more than C. If user specified incoherency bound for the query q is C q, then the dissemination network has to ensure that: () ()1 q q q i i i n i s t d w c $=? \times ??$ (2)

Whenever data values at sources change such that query incoherency bound is violated, the updated value(s) is disseminated to the client. If the network of aggregators can ensure that the i th data item has incoherency bound C i then the following condition ensure that the query incoherency bound C q is satisfied: 1 q q q i i n i c w c = \times ? ?(3)

¹¹⁷ For additive aggregation queries, a client specified query incoherency bound needs to be translated into ¹¹⁸ incoherency bounds for individual data items or sub-queries such that these satisfy Equation (3).

MIN/MAX queries involve of data items, whose extremes are the required result, and its incoherency bound. In a MIN (MAX) query, even if one data value changes it is possible that that value is minimum (maximum) thus individual data incoherency bound cannot be more than query incoherency bound. Thus in case of MIN/MAX queries the dissemination network has to ensure that Ci ? Cq for all the data items appearing in the query.

¹²³ 4 b) Summary of Distributed Execution approach

Consider a client query Q=50 D1 + 200 D2 + 150 D3 with a required incoherency bound of 80 (in a stock portfolio D1, D2, D3 can be different stocks and incoherency bound can be \$80). We want to execute this query over data aggregators given in Example1, minimizing number of refreshes. There are various options for the client to get the data items.

The client may get the data items D1, D2 and D3 separately. The query incoherency bounds can be divided among data items in various ways while satisfying Equation 3. In this report, we show that getting data items independently is a costly option. This strategy ignores facts that the client is interested only in the aggregated value of the data items and various aggregators can disseminate more than one data item.

If a single DA can disseminate all three data items required to answer the client query, the DA can construct 132 a composite data item corresponding to the client query (Sq=50 D1 + 200 D2 + 150 D3) and disseminate the 133 134 result to the client so that the query incoherency bound is not violated. It is obvious that if we get the query result from a single DA, the number of refreshes will be minimum (as in this case data item updates may cancel 135 136 out each other, thereby keeping the query result within the incoherency bound). As different data aggregators 137 disseminate different subsets of data items, no data aggregator may have all the data items required to execute 138 the client query, which is indeed the case in Example1. Further, even if an aggregator can disseminate all the data items, it may not be able to satisfy the query coherency requirements. In such cases, the query has to be 139 executed with data from multiple aggregators. 140

Another option is to divide the query into a number of sub-queries and get their values from individual DAs. In that case, the client query result is obtained by combining the results of more than one sub-query. For the DAs given in Example1, the query Q can be divided in two alternative ways:Plan1: A1 {50 D1 + 150 D3 }; D2 {D2 } Plan2: A1 {D3 }; D2 {50 D1, + 200 D2 }

i.e., in plan1 result of sub-query 50 D1 + 150 D3 is served by A1 whereas value of (or 200 D2) by D2 is
served by A2. In plan2, value of D3 is served by A1 whereas result of sub-query 50 D1 + 200 D2 is served by
A2. Combining them at the client gives the query result.

Selecting the optimal plan among various options is not-trivial. As a thumb-rule, we should be selecting the 148 plan with lesser number of sub-queries. But that is not guaranteed to be the plan with the least number of 149 messages. Further, we should select the subqueries such that updates to various data items appearing in a sub-150 query have more chances of cancelling each other as that will reduce the need for refresh to the client (Equation 151 2). In the above example, if updates to D1 and D3 are such that when D1 increases, D3 decreases, and vice-versa, 152 then selecting plan1 may be beneficial. We give an algorithm to select the query plan based on these observations. 153 While solving the above problem of selecting the optimal plan we ensure that each data item for a client query 154 is disseminated by one and only one data aggregator. Although a query can be divided in such a way that a 155 single data item is served by multiple DAs multiple aggregators, increasing the unnecessary processing load. By 156 157 dividing the client query into disjoint sub-queries we ensure that a data item update is processed only once for 158 each query (For example, in case of paid data subscriptions it is not prudent to get the same data item from the 159 multiple sources).

The query incoherency bound needs to be divided among sub-query incoherency bounds such that, besides satisfying the client coherency requirements, the chosen DA (where the sub-query is to be executed) is capable of satisfying the allocated subquery incoherency bound. For example, in plan1 allocated incoherency bound to the sub-query 50D1 + 150D3 should be greater than 55 (=50*0.5+150*0.2) as that is the tightest incoherency bound which the aggregator D1 can satisfy. We prove that the number of refreshes depends on the division of the query incoherency bounds among sub-query incoherency bounds. Thus, what we need is a method of (a) optimally dividing client query into sub-queries and (b) assigning incoherency bounds to them; such that (c) selected sub-queries can be executed at chosen. And (d) total query execution cost, in terms of number of refreshes, is minimized.

169 **5 II.**

170 6 DATA DISSEMINATION COST MODEL

171 Cost of disseminating a data item at a certain given incoherency bound C can be estimated by combining two 172 models:

¹⁷³ 7 a) Incoherency bound model

The incoherency bound model is used for estimating dependency of data dissemination cost over the desired incoherency bound. As per this model, we have shown in [13] that the number of data refreshes is inversely proportional to the square of the incoherency bound (1/C2). Similar result was earlier reported in [4] where the data dynamics was modeled as a random walk process.Data dissemination cost ? 1/C 2 (4) b) Data synopsis Model

The Data synopsis model is used for estimating the effect of data dynamics on number of data refreshes. We define a data dynamics measure called, sumdiff, to obtain a synopsis of the data for predicting the dissemination cost. The number of update messages for a data item is likely to be higher if the data item changes more in a given time window. Thus we hypothesize that cost of data dissemination for a data item will be proportional to data synopsis, called sumdiff, defined as: () 1 s i i R s s ? = ? ?(5)

Where S i and S i-1 are the sampled values of the data item at i th and (i-1) th time instances (consecutive ticks). Data sumdiff can be maintained at the source or aggregators. For calculating this quantity, the data source can accumulate the absolute value of changes in data items or the data aggregator can estimate this quantity using changes in pushed values. Next we use this result for developing the query cost model.

Consider a case where a query consists of two data items P and Q with weights w p and w q respectively; and we want to estimate its dissemination cost. If data items are disseminated separately query sumdiff will be:

Instead, if the aggregator uses the information Instead, if the aggregator uses the information that client is interested in a query over P and Q (rather than their individual values), it makes a composite data item w pP + w q q and disseminates that data item then the query sumdiff will be: (7) R query is clearly less than or equal compared to R data. Thus we need to estimate the sumdiff of an aggregation query (i.e., R query) given the sumdiff values of individual data items (i.e., R p and R q). Only data aggregators are in position to calculate R query as different data items may be from different sources.

196 **8 III.**

197 9 QUERY COST MODEL

For getting an estimation of the query dissemination cost what we need is R query whereas we know R p and R q (in Equation (6) and (7). As different data items may be disseminated by different servers, R query can be calculated only at data aggregators. If two data items are correlated such that if value of one data item increases other also increases, then R query will be closer R data ; whereas if the data items are inversely correlated then R query will be much less than R data . Thus, institutively, we can represent the relationship between R query and sumdiff of individual data items involved using a correlation measure between data items. Specifically, if ? is the correlation measure then R query can be written as:

205 ()2 2 2 2 2 2 2 p p q q p p q q query R w R w R w R w R ?? + +(8)

The correlation measure is defined such that-1??? +1, so, Rquery will always be less than $|w| p \ge p + w = q$ 206 R q | (as explained earlier) and always be more than |w p R p ? w q R q |. R data = w p R p + w q R q 207 =w p ?| pi -p i-1 | +w q ?| q i ?q i-1 | R query =?|w p (p i ?p i-1)+w q (q i ?q i-1)| by data items P and 208 Q. Cosine similarity is a widely used measure in information retrieval domain where documents are represented 209 using a vector-space model and document similarity is measured using cosine of angle between two document 210 211 representations. For data streams P and Q, ?? can be calculated as: (9) a) Executing queries using sub queries 212 For executing an incoherency bounded continuous query, a query plan is required which includes the set of sub-213 queries, their individual incoherency bounds and data aggregators which can execute these sub-queries. We need 214 to find the optimal query execution plan which satisfies client coherency requirement with the least number of refreshes. What we need is a mechanism to: While satisfying the following conditions identified in Section 1.2: 215 Condition 1. Query incoherency bound is satisfied. 216

Condition 2. The chosen DA should be able to provide all the data items appearing in the sub-query assigned to it. Condition 3. Data incoherency bounds at the chosen DA should be such that the sub-query incoherency bound can be satisfied at the chosen DA. Objective: Number of refreshes should be minimized.

²²⁰ 10 b) Minimum Cost

Figure 2 shows the outline of greedy heuristics where different criteria (?) can be used to select subqueries. In

this section we describe the case where the estimate of query execution cost is minimized in each step of the algorithm (min-cost) whereas in the next section we present the case where gain due to executing a query using sub-queries is maximized (maxgain).

²²⁵ 11 c) Query Plan with Pre-decided Incoherency Bound Alloca tion

For the given client query (q) and mapping between data aggregators and the corresponding {data-item, data incoherency bound } pairs (f : D?(S, C)) maximal sub-queries can be obtained for each data aggregator. Let A be the set of such maximal sub queries. In this set, each query a ? A can be disseminated by a designated data aggregator at the assigned incoherency bound. For each sub-query a ?A, its sumdiff Ra is calculated. Using the set A and sub-query sumdiffs, we use the algorithm outlined in Figure 2 to get the set of sub-queries minimizing the query cost.

²³³ 12 In this Figure each sub-query a ? A is represented by

the set of data items covered by it. As we need to minimize the query cost, a sub-query with minimum cost per data item is chosen in each iteration of the algorithm i.e., criteria ? ? minimize (Ra/Ca 2 |a|).

All data items covered by the selected sub query are removed from all the remaining sub-queries in A before performing the next iteration.

238 13 Algorithm:

Result The decision is taken based on client query information. The greedy method is the most straight forward 239 method. It is popular for obtaining the optimized solutions. In the greedy method there are some important 240 activities. (a) A selection of solution from the given input domain is performed. (b). The feasibility of the 241 solution is performed and then all the feasible solutions are obtained. (c) From the set of feasible solutions, the 242 particular solution that minimizes or maximizes the given objective function is obtained. Such a solution is called 243 optimal solution. For an algorithm that uses greedy method works in stages. At each stage only one input is 244 considered at a time. Based on this input it is decided whether particular input gives the optimal solution or 245 not. 246

²⁴⁷ 14 d) Maximum Gain

In this section we present an algorithm which, instead of minimizing the estimated query execution cost, maximizes the estimated gains of executing client query using sub-queries. In this algorithm, for each subquery, we calculate the relative gain of executing it by finding the sumdiff difference between cases when each data item is obtained separately and when all the data items are aggregated as a single sub-query. (i.e., maximal sub-query).

²⁵³ 15 IV.

²⁵⁴ 16 Related work

²⁵⁸ 17 Answering Incoherency bounded aggregation queries

Various mechanism for efficiently answering incoherency bounded aggregation queries over continuously changing data items are proposed in the literature [3,4,9]. i.e., in this thesis, to develop and evaluate client-pull-based techniques for refreshing data so that the results of the queries over distributed data can be correctly reported, conforming to the limited incoherency acceptable to the users. Here considered the problem of answering queries for online decision making at web data aggregators.

264 Our work distinguishes itself by employing subquery based evaluation to minimize number of refreshes. 265 Pull based data dissemination techniques, where client or data aggregators pull data items such that query 266 requirements are met, are described in [3]. For minimizing the number of pulls, both predict data values and pull instances. In comparison, we use push based mechanism to refresh sub-query values at the client. In [4], 267 268 authors propose push based scheme using data filters at the sources. i.e., distributed data sources continuously stream updates to a centralized processor that monitors continuous queries over the distributed data. Based on 269 we specified a new approach for reducing communication cost in an environment of centralized continuous query 270 processing over distributed data streams. 271

18 B) CONSTRUCTION AND MAINTENANCE OF NETWORK OF DATA AGGREGATORS

This approach hinges on specifying precision constraints for continuous queries, which are used to generate 272 adaptive filters at remote data sources that significantly reduce update stream rates while still guaranteeing 273 sufficient precision of query results at all times. And enables users or applications to trade precision for lower 274 communication cost at a fine granularity by individually adjusting precision constraints of continuous queries. 275 276 Imprecision of query results is bounded numerically so applications need not deal with any uncertainty. To validate our approach we performed a number of experiments using simulations and a real network monitoring 277 implementation approach in achieving low communication overhead. According to that work can an aggregation 278 query, the number of refresh messages can be minimized by performing incoherency bound allocation to individual 279 data items such that the number of messages from different data sources is the same. Instead we execute more 280 281 dynamic assigning incoherency bounds. And minimizing the total number of messages send by DAs. Like us ,authors of [9],also assume that dissemination tree from sensor node[data source] to root[client]already exist; and 282 they also install error filters on partial aggregates (similar to in coherency bound assign to sub queries) but, in 283 our work each data aggregator can only discriminates data at some pre-specified incoherency bound depending 284 on its capability where as such a constraints does not exist for [9].further, we also be give method to select partial 285 aggregates (sub queries) to be used to answering the query. 286

Authors propose using data filters at the sources; instead we assign incoherency bounds to subqueries which reduce the number of refreshes for query evaluation, Further, we propose that more dynamic data items should be executed as part of larger sub-query. In [8], i.e., here discuss various techniques of reorganizing a data dissemination network when client requirements change. Instead, we try to answer the client query using the existing network. Reorganizing aggregators is a longer term activity whereas query planning can be done for short as well as long running queries on more dynamic basis.

Like us, author of [9] also assume that dissemination tree from sensor nodes (data-sources) to root (clients) 293 already exists. In-network data aggregation has been recently proposed as an effective means to reduce the 294 number of messages exchanged in wireless sensor networks. Nodes of the network form an aggregation tree, 295 in which parent nodes aggregate the values received from their children and propagate the result to their own 296 parents. However, this schema provides little flexibility for the end-user to control the operation of the nodes 297 in a data sensitive manner. For large sensor networks with severe energy constraints, the reduction (in the 298 number of messages exchanged) obtained through the aggregation tree might not be sufficient. In this thesis 299 we present new algorithms for obtaining approximate aggregate statistics from large sensor networks. The user 300 specifies the maximum error that he is willing to tolerate and, in turn, our algorithms program the nodes in a 301 way that seeks to minimize the number of messages exchanged in the network, while always guaranteeing that 302 303 the produced estimate lies within the specified error from the exact answer. And they also install error filters on partial aggregates. But in our work, each data aggregators can only disseminate data some pre-specified 304 incoherency bound depending on its capability whereas such a constraint does not exist for [9]. Further, we also 305 give a method to select partial aggregates (sub queries) to be used for answering the query. In [12] Pull based 306 data dissemination techniques, where clients or data aggregators pull data items such that query requirements 307 are met, are described in [3]. i.e., we develop and evaluate client-pull-based techniques for refreshing data so that 308 the results of the queries over distributed data can be correctly reported, conforming to the limited incoherency 309 acceptable to the users. For minimizing the number of pulls, both model the individual data items and predict 310 data values. In comparison, we consider the situation where different sub-queries, involving multiple data items, 311 can be evaluated at different nodes. Further, incoherency bound is applied over the sub-query rather than to 312 Spatial and temporal correlations between sensor data are used to reduce data refresh instances in [5,6]. We also 313 consider correlation in terms of cosine similarity between data items, but we use it for dividing client query into 314 sub-queries. 315

³¹⁶ 18 b) Construction and maintenance of network of data aggre ³¹⁷ gators

Authors of [1,2,8] describe Construction and maintenance of hierarchical network of data aggregators for providing scalability and fidelity in disseminating dynamic data items to large number of clients. In these works, fidelity is defined as fraction of time when the client coherence requirements are met. Each data aggregators is given client requirements in the form of data items and their respective incoherency bounds. Instead we use such networks for efficiently answering client's aggregation queries.

One can use client queries to optimally construct a network of data aggregators while, on the other hand, one can also use a given network of aggregators to efficiently answer client queries. Authors of [1,2,8] deal with the first part where as we have studied the second part. Changes in data dynamics may lead to reorganization of the network of data aggregators which, in turn necessitate changes in query plans. Whereas query plan can change more often depending on data dynamics.

Instead of optimizing fidelity of data items at data aggregators, as proposed in [2], using our work, one can optimize fidelity all the way up to client queries. Fidelity of a data item can be approximately calculated as number of dissemination messages multiplied by the total delay in the message transmission. Author of [2] assume that each client's requirements are fulfilled by a single data aggregator. But in case of data aggregators may need to disseminate a large number of data items which will lead to processing large number of refresh messages,

hence increase in delay. Thus, each client getting all its data items from a single data aggregators(using a single 333 sub-query) is optimal from number of messages point of view but not necessarily from the query fidelity point 334 of view. By using our work, one can model expected number of messages for client query. Thus, our work can 335 complement the of [2] for end-to-end (source-to-client) fidelity optimization. 336 ν.

337

Conclusion and future work 19 338

In this literature presents a cost based approach to minimize the number of refreshes required to execute an 339 incoherency bounded continuous query. For optimal execution we divide the query into subqueries and evaluate 340 each sub-query at a chosen aggregator. Performance results show that by our method the query can be executed 341 using less than one third the messages required for existing schemes. 342

Further we showed that by executing queries such that more dynamic data items are part of a larger sub-query 343 we can improve performance. Our query cost model can also be used for other purposes such as load balancing 344 various aggregators, optimal query execution plan at an aggregator node, etc. 345

Developing efficient strategies for multiple invocations of our algorithm, considering hierarchy of data 346 aggregators. Another area for future research is changing a query plan as data dynamics changes. Another 347 area of our future work is using the cost model for these applications and developing the cost model for more 348 complex queries.



Figure 1:

- 1350 [NEFSC Scientific Computer System], http://sole.wh.whoi.edu/jmanning//cruise/servelet.cgi 1351 NEFSC Scientific Computer System
- 352
 [Pearson Product moment correlation coefficient] ,
 http://www.nyx.net/~tmacfarl/STAT_TUT/

 353
 correlat.ssi/ Pearson Product moment correlation coefficient
- [Fei ()] A Novel Approach to Managing Consistency in Content Distribution, Zongming Fei . 2001. WCW.
- [Olston et al. ()] Adaptive Filter for Continuous Queries over Distributed Data Streams, C Olston, J Jiang, J
 Widom. 2003.
- ³⁵⁷ [Hochbaum ()] 'Approximation algorithms for the set covering and vertex cover problems'. D S Hochbaum .
 ³⁵⁸ SIAM Journal on Computing 1982. 11 (3) .
- [Shah et al. ()] Client Assignment in Content Dissemination Networks for Dynamic Data, S Shah , K Ramam ritham , C Ravishankar . 2005.
- 361 [Agrawal et al. ()] Construction of a Temporal Coherency Preserving Dynamic Data Dissemination networks, S
- Agrawal, K Ramamritham, S Shah. 2004. (RTSS)
 [Davis et al. ()] Edge Computing: Extending Enterprise Applications to the Edge of the Internet, A Davis, J
- Parikh , W Weihl . 2004. WWW.
 [Gupta et al. ()] Executing Incoherency Bounded Continuous Queries at Web Data Aggregators, R Gupta , A
- [Gupta et al. ()] Executing Incoherency Bounded Continuous Queries at Web Data Aggregators, R Gupta , A
 Puri , K Ramamritham . 2005. WWW.
- ³⁶⁷ [Dilley et al. (2002)] 'Globally Distributed Content Delivery'. J Dilley , B Maggs , J Parikh , H Prokop , R
 ³⁶⁸ Sitaraman , B Weihl . *IEEE Internet Computing* Sept 2002.
- [Shah et al. ()] Maintaining Coherency of Dynamic Data in Cooperating Repositories, S Shah , K Ramamritham
 , P Shenoy . 2002.
- ³⁷¹ [Optimized Execution of Continuous Queries, APS2006] Optimized Execution of Continuous Queries, APS2006,
 ³⁷² www.cse.iitb.ac.in/~grajeev/APS06.PDF
- [Guptha and Ramamritham ()] Optimized Query Planning of Continuous Aggregation Queries in Dynamic Data
 Dissemination Networks, R Guptha , K Ramamritham . 2007. WWW.
- [Vander Meer et al. (2004)] 'Proxy-Based Acceleration of Dynamically Generated Content on the World Wide
 Web'. D Vander Meer , A Datta , K Dutta , H Thomas , K Ramamritham . ACM Transactions on Database
 Systems (TODS June 2004. 29.
- 378 [Query cost model validation for sensor data] Query cost model validation for sensor data, www.cse.iitb.ac. 379 in/~ravivj/BTP06.pdf
- [Guptha and Ramamritham] Query Planning for Continuous Aggregation Queries over a network of Data
 Aggregators, R Guptha, K Ramamritham. IEEE 2011.
- 382 [Rangarajan et al. ()] 'User Specific Request Redirection in a Content Delivery Year Network, 8th Intl'. S
- Rangarajan , S Mukerjee , P Rodriguez . Workshop on Web Content Caching and Distribution (IWCW),
 2003.