# Regional Digest Aggregation based on Opportunities in Wireless Sensor Networks

**Chi Quynh Nguyen, Ngoc thi Bich Do**

***Abstract**: It is desirable to reduce the amount of data collected in sensor networks to reduce energy consumption and to extend network lifetime. At the same time, one should extract as much information as possible to allow support for heterogeneous user queries. Several aggregation proposals have been proposed to reduce the amount of data communication within sensor networks, primarily focusing on supporting limited and straightforward query types, such as SUM, COUNT, AVG, and MIN/MAX. Unfortunately, user queries are not limited to these simple types of aggregates and cannot be predicted a priori. In this paper, we propose an aggregation framework that produces regional digests at a parameter-defined granularity, allowing for the support of arbitrary user queries. Since the success of the aggregation policy greatly depends on the integrated routing mechanisms, we evaluate the performance of our approach under alternative routing approaches. Our experimental results suggest at least 3-fold improvement in spatial accuracy at a relatively small expense of increased energy consumption.*

***Keywords**: Sensor Networks, In-Network Data Management, Aggregation, Adaptive Routing, Energy-Efficiency.*

## I. INTRODUCTION

Due to advances in sensing equipment, sensor networks are becoming a highly popular tool for various scientific, commercial, and military applications. Environmental monitoring, target tracking, disaster monitoring, and earthquake and structural engineering are just a few examples of such applications. Regardless of the particular application, the main objective of sensor networks is to forward observations made in an area to a super node to which users can issue queries. When data collection is initiated with an explicit request before data collection, it is referred to as *pull*, whereas when sensor nodes forward their observations without an explicit request, it is referred to as *push* [1]. The success of any sensor network depends on its ability to respond to user queries on the data collected using either a pull or push-based approach. In this regard, it is preferable to extract as much helpful information as possible from the network, since user queries cannot be predicted a priori.

*\*Correspondence Author(s)*
**Chi Quynh Nguyen**\*, Faculty, Department of Computer Science, Posts and Telecommunications Institute of Technology, Hanoi, Vietnam. Email: chinq@ptit.edu.vn, ORCID ID: 0009-0007-6197-2486
**Ngoc thi Bich Do**, Faculty of Department, Posts and Telecommunications Institute of Technology, Hanoi, Vietnam. Email: ngocdtb@ptit.edu.vn, ORCID ID: 0009-0004-3250-0154

To date, a significant focus in sensor network studies has been to minimise energy consumption in the network, as ordinary nodes responsible for making observations are typically powered by batteries. Batteries have a limited lifespan and may often be irreplaceable or non-rechargeable. Therefore, for this type of power source, it is preferable to slow down the rate at which power is drained from batteries to extend the network lifetime. As transmitting one bit over radio is three orders of magnitude more expensive in terms of energy consumption [6], it is desirable to reduce the amount of data traffic during data collection. In this study, we propose an energy-efficient regional digest algorithm, R-Digest, that is integrated with routing for efficient data collection, enabling users to issue a wide range of spatial queries after data collection. Using a tunable aggregation parameter, we can trade off between query response accuracy and energy efficiency. We study the impacts of such aggregation granularity for alternative routing policies. Our performance evaluation results suggest that R-Digest has significant potential to produce reasonable approximations for improved query support. We have observed at least a 3-fold improvement in spatial query accuracy at a very reasonable increase in energy consumption. The rest of the paper is organized as follows. In section 2, we present our motivation behind R-Digest, the aggregation policy we propose. In Section 3, we outline the R-Digest aggregation data structures and the query evaluation process. Section 4 presents the highlights from our performance evaluation for various aggregation parameters and routing policies. Finally, we conclude with Section 5.

## II. RELATED WORK AND MOTIVATION

### A. Existing Proposals on Data Aggregation

The main objective of wireless sensor networks is to make observations in the area of deployment and to communicate these observations to users who access and query the collected data. In general, observations made are collected using either a push-based or a pull-based approach [1]. In the pull-based approach, queries are issued on demand and are directed to a specific geographic area according to selected information sources [2][3]. In the case of a push-based approach, on the other hand, data collection is initiated without an explicit request from the clients. Therefore, the types of queries that need to be supported are not known a priori. In this respect, the original data should be maintained as much as possible, so that arbitrary queries that may arrive later can be supported. In the following, we discuss representative issues related to data access.

# Regional Digest Aggregation based on Opportunities in Wireless Sensor Networks
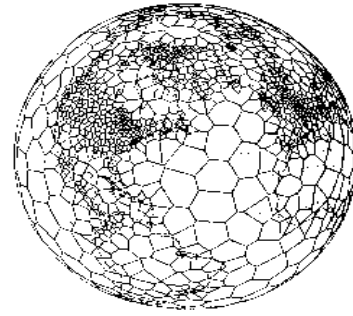
Data aggregation, which combines messages from multiple nodes along the routing path, is a popular tool for reducing the amount of data that needs to be communicated. To date, several attempts have been made to utilise data aggregation. Routing protocols such as LEACH [4] or PEGASIS [5] are well-known examples that have considered the possibility of such aggregation. However, a large number of data aggregation policies assume that all data collected is perfectly compatible and can be aggregated without any loss. For instance, Lindsey et al. [5] have assumed that any two messages received at any node can be combined into a single message. In this study, we refer to such aggregates as *singleton aggregates*, as they aim to produce a single aggregate value for the entire network. Q-Digest [6] proposes an attempt to acquire more detailed information from the sensor network by utilising histogram-based digests within the network. However, histogram-based digests are focused on a particular aggregation function and are not sufficient to respond to arbitrary queries. Synopsis Diffusion [7] is a recent study that defines a framework for aggregates, yet it focuses on simple functions such as SUM, COUNT, AVG, MIN/MAX. A thorough survey on the basic techniques for data aggregation can be found in [9]. Recently, researchers have become more concerned with methods for scheduling the deployment of the data aggregation process to minimise query processing time [10][11].

Although spatial correlation has been considered in some work, such as [12], in essence, none of existing solutions provide sufficient detail to respond to spatial (location dependent) queries which are the most common query types for expert applications. Examples can be in the form "*What was the pressure measured in the region that has reported the lowest temperature?*", or "*What is the chemical dispersion rate at the southeast quadrant of the lake?*", etc.
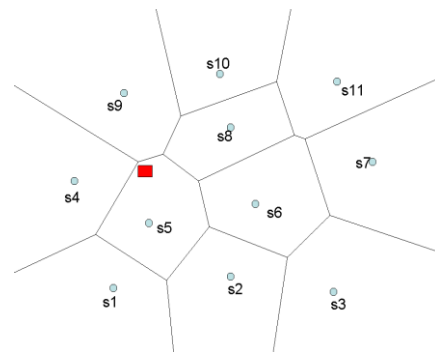
## B. Motivation: Responding to Spatial Queries

For a moment, let us assume that there are no resource limitations in the network and *all* of the observations can be reported to data collection points. Even in cases where all individual observations are available at data collection points, query responses may still require some processing. In a sensor network application, only discrete points can be measured, while queries can be directed to any point in the region. In such cases, we need to extrapolate the measured readings to other geographic areas. In this regard, one approach used is based on Voronoi diagrams.

In fig. 1, an example scenario is depicted where 11 sensors are deployed in a field represented by small circles. These represent a superset of the discrete points from which we can obtain an observation in the field. The small rectangle represents a point of interest for a user query. As demonstrated, the user queries are not limited to the discrete data points collected from the field. In this regard, different approximation policies are possible. In the figure, each sensor reading is assumed to apply to all points that are closest to its location in comparison to any other sensor location. In other words, for any point of interest for a particular query, we assume that the closest reading applies.



**(a) Voronoi Diagram of the Earth, Dividing the Surface into Cells**



**(b) A Sample Scenario Using 11 Sensors Deployed in the Field**

**Fig.1. Observations collected by individual sensor nodes are only discrete samples of a continuous space. User queries, on the other hand, can be issued for arbitrary locations in the network. This requires an approximation based on discrete measurements.**

In this approach, the correlation between neighbouring nodes is entirely ignored. For instance, in Figure 1, assume that the readings at sensor nodes S4, S5, and S9 are 11, 6, and 12, respectively. Using the Voronoi graph approach, the query will be replied to as 6, ignoring the fact that the two other sensors nearby have significantly higher readings.

## C. Implications for Large-Scale Deployments

Large-scale deployments of sensor networks involve a vast number of sensors that span a broad geographic area. As an example, PLASMA (PLAnetary Scale Monitoring Architecture) [2] is an interdisciplinary project that aims at an integrated architecture of heterogeneous sensor networks in highly distributed environments. In large-scale deployments, even if it were possible to obtain all observations, it is still not preferable to do so due to the enormous amount of data that would exceed the servers' processing capacity. Therefore, we try to push data processing closer to the source, the sensor node itself. This would not only reduce energy consumption in the network but also help optimise query processing costs at the data collection points. As a result, a significant amount of approximation is expected when collecting data to respond to individual queries.

Yet it is not possible to accurately predict user queries that will arrive later. In short, it is not preferable to trade off data expressiveness for energy optimisations. Previous aggregation approaches aim to optimise energy consumption in the network at the expense of only being able to reply to limited queries. For large-scale deployments of sensor networks, we need to pay particular attention to the trade-off between efficiency and query response accuracy.

### III. R-DIGEST

The primary objective of R-Digest is to enable spatial content in the aggregates being created, allowing for the support of arbitrary queries. For this purpose, we first analysed the most descriptive and efficient way to respond to point queries targeted to a small area, as demonstrated in Fig. 1. Due to the inherent ambiguity in combining data from multiple sources, we employ an opportunistic approach that exploits the correlated nature of sensor readings.
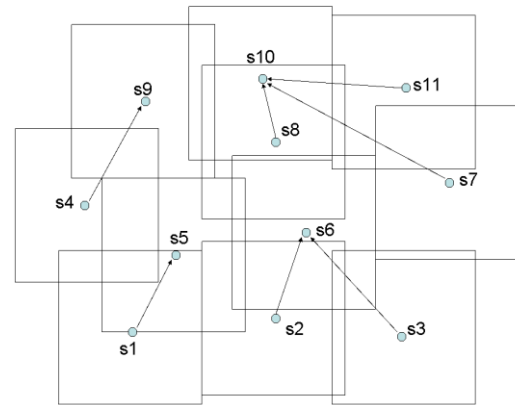
In R-Digest, we define a range of impact around each sensor node during the network setup phase. For simplification of the data structures and thereby the ordinary node processing, we define this region as a perfect square, which is a bounding box around the sensor node. [1] . Alternative definitions are, of course, possible provided that ordinary sensor processing is efficient. A bounding box is defined by a set of two points, the low, i.e., southeast, and the high, i.e., northwest, point. The size of the bounding box has a default value of *the area covered divided by the number of sensors* and can be adjusted according to the certainty level of a reading.
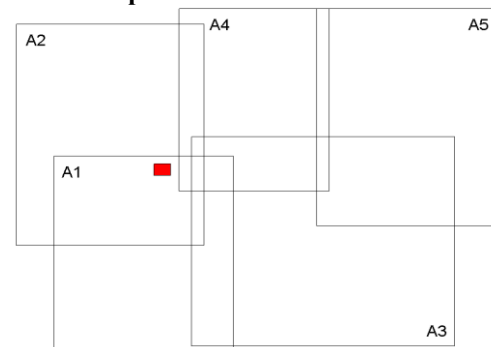
Each node reports an observation in the form:

*<value, min_X, min_Y, max_X, max_Y, first_time, last_time, cnt >*

Where *value* is the specific reading being reported, (*min_X, min_Y)* is the lowest and (*max_X, max_Y)* is the highest coordinate within the bounding box, *first_time* is the time stamp of the earliest observation, *last_time* is the time stamp of the latest observation represented in the record (initially *first_time* and *last_time* are the same), and *cnt* is the total number of nodes that contribute to this reporting (initially 1). It is possible to extend the data structure to include additional statistics, e.g., the minimum, the maximum, etc., readings in the field[2]. 

In fig. 2, we revisit the example scenario with 11 sensors deployed in the field. Assume that the initial bounding box size is determined to be four, and node s1 is located at coordinate (5,6). This sensor will then report its reading of 9 at time 2, with the record being < 9, 3, 4, 7, 8, 2, 2, 1>. In the figure, we annotate the figure with current data communications taking place. An edge between two nodes suggests that the data record of the node is forwarded to the node to which it is connected. For instance, node s1 will forward its record to node s5.



**Fig. 2. Each Node Reports its Observation Using A Bounding Box Using Homogeneous Box Sizes. In the Figure, A Snapshot of the Current State of Routing Around is Also Depicted**



**Fig. 3. Data Aggregation, Using A Parameter of 2, Results in 5 Aggregate Records, rather than 11 Individual Observations, to be Collected After Observations are Forwarded Using the Depicted Paths in Fig. 1. Spatial Queries Will be Responded Based on the Resulting Aggregates**

During routing, multiple readings will be combined in a single record according to the aggregation parameter. For instance, when node s1 is forwarding its observation to node s5, whose initial reading is <7, 4, 4, 8, 10, 3, 3, 1>, the two readings can be aggregated. For an aggregation parameter of 2, such that the maximum bounding box for aggregates is allowed to be as large as twice the original size, the aggregate record A1 will look like < 8, 3, 4, 7, 10, 2, 3, 2 > suggesting that two observations are aggregated in the region with an average reading of 8. In Fig. 3, we plot the aggregates generated from 11 individual readings as a result of the routing depicted in Fig. 2. In this example, five aggregate records were generated out of 11 records using an aggregation parameter of 2. As demonstrated in this example, the size of the aggregated message at a particular node is not limited to one record. In particular, node s10 has received three additional records. Yet, it cannot reduce the four observations down to one due to the restriction on the aggregation parameter. The specific number of aggregate records depends on the arrival time of the observation and the routing policy being applied. For instance, for the same topology presented in Fig. 2, it is possible to reduce the number of aggregates down to 4 using alternative routing approaches.

---

[1] We have observed that for skewed deployments heterogeneous box sizes provide a finer granularity in aggregates.

[2] In the rest of the discussion, we limit the definition to the average reading for simplification.

As demonstrated, a significant question in data aggregation is the impact of such aggregation when combined with data routing in the network. In particular, routing policy has a direct effect on the effectiveness of the aggregation. In our experiments, we study the impact of alternative routing policies on data aggregation. In particular, we study a traditional routing protocol that is designed independently of the data collection pattern. We then focus on a routing policy, FHTL (First Hop Then Leap), based on opportunistic routing [8] that takes data generation pattern into account. FHTL adapts routing paths and data transfer according to the data generated.

Data collection using the described integrated routing and aggregation will result in several aggregate and individual records. Once such observations are extracted from the sensor network itself, no further aggregation is applied at the data collection point, as these nodes do not have the same resource constraints as ordinary sensor nodes. At this point, the focus shifts to the flexibility of the data structure for various queries rather than the cost of data collection. In this regard, it should be possible to reconstruct the area covered by sensor nodes to approximate the original network readings.

For this purpose, we maintain records organised by an index that enables arbitrary spatial query processing. For instance, to reconstruct the reading within the marked query region in fig. 3, *we will retrieve the records for A1 and A2* as they include this region, and report the weighted average of the two aggregates. As a result, the query response will be a function of four individual readings, in contrast to a single reading, as shown in Fig. 1. For regions that overlap with multiple areas, the procedure works without additional complexity, based on the weighted coverage of the area. Using alternative spatial data structures, it is possible to provide the user interface with the complete value distribution based on collected data, enabling continuous monitoring of the area.

As can be easily seen in this example, unlike singleton aggregation techniques, in R-Digest, the approximation error in the regional query response is limited to that within the close neighbourhood, rather than the entire network. A significant question we analyze for performance evaluation is the impact of such approximations. In the next section, we present the results from our experiments.
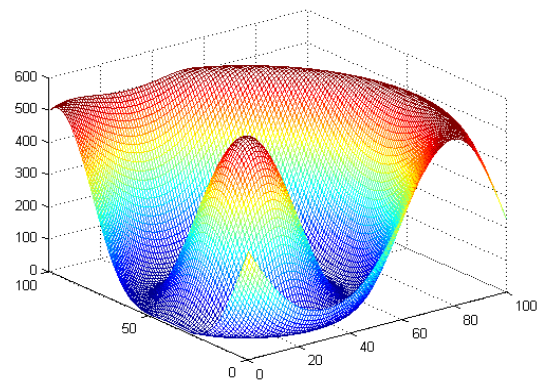
## IV. PERFORMANCE EVALUATION

To study larger-scale networks that our accelerators cannot enable, we have implemented a simulator using the C++ programming language to evaluate the performance of our algorithm. We compared our algorithm to Synopsis Diffusion [7], which is modelled for the perfect average function with unique counting properties. These approaches have their multi-hop routing schemes. A significant question we examine is the impact of routing on data collection and energy efficiency. For this reason, we adapted two approaches for comparison:

As the base case, *R-Digest(L),* we apply LEACH [4] routing and form clusters by electing cluster heads based on a probability value that favours nodes with the most energy. After the selection of cluster heads, nodes join the heads to form clusters. Cluster heads utilise long-range transmission

to communicate directly with the data collection point. Note that alternative descriptions also exist that enforce multi-hop routing from cluster heads. However, these alternatives are less efficient due to the overuse of neighbour communication, which is very costly when aggregation is not possible because of the diversity of the data collected. In our experiments, we exclude the control message overhead to form and dynamically alter the clusters. Therefore, the results presented are only an upper bound on its performance.

As an alternative routing approach, FHTL (First Hop Then Leap), we initiate regular multi-hop routing for several rounds, after which further aggregation becomes unlikely. Beyond this point, we forward the observations directly to the data collection points using direct communication when necessary, without investing in control messages for static organisation. This approach is referred to as *R-Digest (FHTL).*



**Fig. 4. Original Values in the Field Using a Continuous Function**

The main metrics we consider are the accuracy for query support and the overall network lifetime. To evaluate the accuracy of the query responses, we used the following strategy. First, we constructed a correlated value space within the complete network, as illustrated in Fig. 4. This 3D graph represents the ideal case where sensors can be deployed so densely that the continuous value space can be fully defined. In practice, we can deploy only a finite number of sensors in a given area. In Fig. 4, we present a snapshot of the value matrix, which represents the area being monitored. This figure represents the continuous function that the sensor nodes are trying to capture at discrete points.

In the experiments, we used 500 nodes randomly distributed in a 500 m x 500 m grid—the data collection point was located at the southwest corner (0, 0). The initial bounding box size is 10. For the first experiment, we uniformly generate one packet of data at each node every 10 rounds. A round is the time unit we used in our experiments, and each transmission takes one round. At the beginning of the experiment, each node has $2.5 \times 10^6$ nJ of energy.

Each transmission expends an energy of $E_{Tx}(k,d) = E_{elec} * k + \varepsilon_{amp} * k * d^2$ for the transmitting node and $E_{Rx}(k,d) = E_{elec} * k$ for the receiving node. $E_{elec}$ is 50 nJ/bit, and $\varepsilon_{amp}$ is 100 pJ/bit/m$^2$ for the transmit amplifier to achieve an acceptable signal-to-noise ratio.

When considering the reconstructed value space based on the aggregated values at the data collection point, a small aggregation parameter yields a very close approximation. For larger values of the aggregation parameter, the approximation level increases. As a result, a significant question we investigated in our performance evaluation is the trade-offs of the parameter setting.

Our primary performance metric is the accuracy of query responses based on the collected data. In fig. 5, we plot the maximum deviation from the actual values as the aggregation parameter of R-Digest is increased. As expected, the accuracy degrades as this parameter is increased. However, this is at a relatively slow rate, and the overall benefits are at least 3-fold within the whole range in comparison to the alternative, Synopsis Diffusion. The global aggregate approaches, such as SUM and MAX, result in significant deviations for regional queries.

A more interesting observation is that the error in regional queries for R-Digest is only slightly affected by alternative routing policies. Interestingly, the traditional aggregation policy of LEACH favours regional aggregation, which in turn improves accuracy for spatial queries. The more sophisticated routing policy of FHTL, although it significantly improves energy savings, as will be analysed next, introduces an additional approximation for neighbouring nodes that are part of separate routing paths.
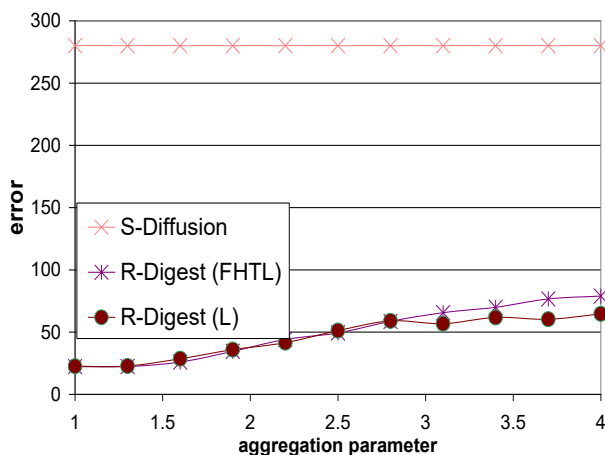


**Fig. 5. Deterioration in Accuracy as the Overall Area is Reconstructed**

To analyse the impacts of energy consumption, we focused on applications that require a wide coverage of the monitored area. Traditional energy consumption evaluation approaches focus on the number of nodes still alive in the network. For this purpose, an implicit threshold is used to describe the end of the network's useful lifetime. For instance, when 70% of the nodes have depleted their energy, the network is considered dead.

In PLASMA, we found that the location of the live nodes is as essential as the number of such nodes. In particular, having 30% of nodes within 5% of the total area is not as useful as having 30% of nodes alive scattered in a larger area. Based on this observation, we focus on the area covered by live nodes as it defines our remaining monitoring capacity within the deployment area.
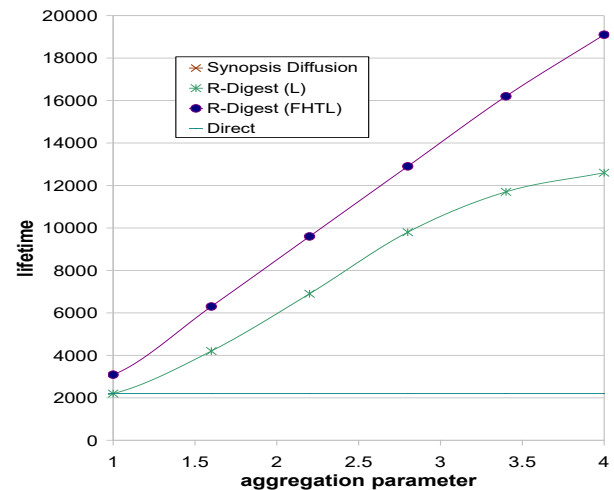


**Fig. 6. Network Lifetime Due to the Distribution of Energy Consumption in the Network**

In fig. 6, we plot the network lifetime as the aggregation parameter is increased. We define network lifetime according to the distribution of nodes in the network, rather than a scalar value representing the number of live nodes. This approach enables us to accurately reflect the distribution of live nodes in the network, allowing for the representation of unbalanced energy consumption schemes.

For this metric, we measure the percentage of the area still covered by at least one live node within a radius of 10 and report the time when such coverage falls below 40% of the original deployment area.

In these experiments, we focus on the combined impact of the routing policy and the particular aggregation technique being employed. In Fig. 6, we observe that R-Digest, when used with LEACH routing, starts to level off when the aggregation parameter is increased beyond 3. On the other hand, the lifetime increases linearly for the FHTL routing model. As demonstrated for larger values of the aggregation parameter, the benefits of R-Digest (FHTL) will continue to increase linearly. When considering the energy consumption of Synopsis Diffusion (indicated by the line with a lifetime of 12000 in Fig. 6), we observe that for small values of the aggregate parameter, the difference is quite significant. Yet the difference collapses when the aggregation parameter is increased. Based on our results, we conclude that R-Digest is a more powerful aggregation policy, allowing for energy savings in direct proportion to the approximation applied, in terms of the data collected. As energy constraints cease to be a limiting factor in sensor networks based on renewable energy sources, such as solar cells, the accuracy of the collected data can have an even more significant impact on performance.

## V. CONCLUSIONS

Data aggregation is a popular approach to reduce the amount of data traffic in the network, thereby improving the network's lifetime. There are two categories of data access for sensor networks: queries on existing or future observations, and those on past observations.

For the second class of queries, the nature of user queries cannot be predicted a priori. As a result, it is desired to maintain as detailed information as possible from sensor observations. In this regard, it is challenging to formulate a successful aggregation function during data collection. In this paper, we propose a new aggregation approach for query processing, called R-Digest, specifically designed for sensor networks. R-Digest produces regional digests during push data collection, allowing for the complete data space to be reconstructed, with some level of approximation, for spatial queries. This approach provides a significant improvement in terms of queries that can be supported using collected data.

## DECLARATION STATEMENT

| Funding | No, I did not receive. |
|---|---|
| Conflicts of Interest | No conflicts of interest to the best of our knowledge. |
| Ethical Approval and Consent to Participate | No, the article does not require ethical approval or consent to participate, as it presents evidence. |
| Availability of Data and Materials | Not relevant. |
| Authors Contributions | All authors have equal participation in this article. |

## REFERENCES

1. Aksoy D., and Leung M.S. - Pull vs Push: A Quantitative Comparison for Data Broadcast. IEEE Global Telecommunications Conference (2004) Volume 3 1464-1468
2. Aksoy D.- PLASMA: A Planetary Scale Monitoring Architecture. Proceedings of ACM international conference on Multimedia (2005) 96-102 https://doi.org/10.1145/1101149.1101164
3. Govindan R, Estrin D, Intanagonwiwat C- Directed Diffusion: A Scalable and Robust Communication Paradigm for Sensor Networks. Proceedings of the international conference on Mobile computing and networking (2000) 56-67. https://doi.org/10.1145/345910.345920
4. Heinzelman WR, Chandrakasan A, Balakrishnan H Energy-Efficient Communication Protocol for Wireless Microsensor Networks, IEEE Hawai Int. Conf. on System Sciences (2000).
5. Lindsey S, Raghavendra CS- PEGASIS, Power Efficient Gathering in Sensor Information Systems, IEEE Aerospace Conference Proceedings (2002) 1125-1130 Vol 3.
6. Shrivastava N., Buragohain C., Agrawal D., Suri S.- Medians and Beyond: New Aggregation Techniques for Sensor Networks. Proceedings of the international conference on Embedded networked sensor systems (2004) 239-249. https://doi.org/10.1145/1031495.1031524
7. Suman Nath, Phillip B. Gibbons, Srinivasan Seshan, and Zachary Anderson - Synopsis Diffusion for Robust Aggregation in Sensor Networks. ACM Transactions on Sensor Networks (2008). Volume 4, Issue 2, Article No. 7. https://doi.org/10.1145/1340771.1340773
8. Chen C., Aksoy D., and Demir T.- Processed Data Collection using Opportunistic Routing in Location-Aware Wireless Sensor Networks. Proceedings of the International Conference on Mobile Data Management (2006) 150-158.
9. Fasolo E., Rossi M., Widmer J., Zorzi M., In-network aggregation techniques for wireless sensor networks: a survey. IEEE Wireless Communication. Vol 14, Issue 2, (2007) 70-87. https://doi.org/10.1109/MWC.2007.358967
10. Yu B., Li J., and **Li Y.**, Distributed Data Aggregation Scheduling in Wireless Sensor Networks, IEEE INFOCOM 2009,(2009) 19-25. https://doi.org/10.1109/INFCOM.2009.5062140
11. Xiaohua Xu, Xiang-Yang Li, Xufei Mao, Shaojie Tang, Shiguang Wang - A Delay-Efficient Algorithm for Data Aggregation in Multihop Wireless Sensor Networks. IEEE Transactions on Parallel and Distributed Systems, Vol. 22, No. 1 (2011) 163- 175 https://doi.org/10.1109/TPDS.2010.80
12. Yujie Zhua, Ramanuja Vedanthama, Seung-Jong Parkb. Raghupathy Sivakumara - A scalable correlation-aware aggregation strategy for wireless sensor networks. Information Fusion, Volume 9, Issue 3, (2008) 354-369. https://doi.org/10.1016/j.inffus.2006.09.002

## AUTHORS PROFILE

**Chi Quynh Nguyen** graduated with a Bachelor of Science in Computer Science from Hanoi University of Technology in Vietnam in 1999, earning a summa cum laude distinction. She then received a Vietnamese Government Fellowship to pursue a Master of Science in Computer Science at the University of California, Davis, USA, in 2004. Then she became a Ph.D. candidate in Computer Science at the same University in 2006. Since 2008, she has been a senior lecturer in the Faculty of Information Technology at the Posts and Telecommunications Institute of Technology in Hanoi, Vietnam. Her primary research focuses on data warehousing, data mining, and bioinformatics, as well as Mobility prediction, self-configuration of MANets, and data aggregation methods in sensor networks.

**Ngoc Thi Bich Do** earned her Bachelor of Science degree in Information Technology from the University of Science and Technology in 2004. Subsequently, she successfully earned her Master's degree in Computer Science from the University of Hanoi in 2007. In 2010, she successfully earned her Ph.D. degree from the Japan Advanced Institute of Science and Technology, specialising in the field of Information Science. Since 2013, she has been an esteemed lecturer within the Faculty of Information Technology at the Posts and Telecommunications Institute of Technology. Her research interests include software testing, formal methods, numerical analysis, data mining, and machine learning.