

INSPIRE-DB: Intelligent Networks Sensor Processing of Information using Resilient Encoded-Hash DataBase

Vasanth Iyer*, S. Sitharama Iyengar[†], Garmiela Rama Murthy*,
Kannan Srinathan*, Vir Phoha[‡] and Mandalika B. Srinivas[§]

*International Institute of Information Technology, Hyderabad, India - 500 032

[†]Louisiana State University, Baton Rouge, LA 70803, USA

[§]Brila Institute of Technology & Science, Hyderabad Campus, Hyderabad-500078, India

[‡]Louisiana Tech University, Center for Secure Cyberspace, LA 70803, USA

vasanth@research.iiit.ac.in; iyengar@csc.lsu.edu; phoha@coes.latech.edu; srinivas@bits-hyderabad.ac.in
{srinathan, rammurthy}@iiit.ac.in

Abstract—Sensor networks consist of small motes attached with sensors to measure ambient parameters like temperature, humidity and light. As these motes are unreliable due to wireless link quality and also the data measuring sensors cannot be calibrated accurately for a given applications need. The unique data fusion needs are that parameter being measured is distributed across the network and needs to be computed reliably and with minimum overhead and redundancy due to data value being correlated. We show the asymptotic complexity of topology control when applied to power-aware routing is scalable and argue that the accuracy and reliability of the estimated sensor values can be accurately predicted for the physical value being sensed and aggregating. A prefix-based routing protocol is used for data-centric storage, which allows querying distributed parameters using a KEY, VALUE pairs without the need of the sensor node to know its exact geographic information. Intelligent sensor information processing, which is driven by these requirements, is discussed under the framework INSPIRE-DB.

Keywords-Data-centric routing; Distributed Hash Table (DHT); Distributed Source Coding (DSC); Distributed Compressed Sensing (DCS); Sensor Fusion; Pre- and Post processing of Sensors; Cross-layer Protocols.

I. INTRODUCTION

The IEEE 802.15.4 specification [8] defines wireless sensor network in terms of low data rate of (100kbps), short radio range of 50 m² and low power consumption of 60–80 (milliwatts) per transmit burst for longest lifetime performance for up to 2³² node topology. Many of the existing routing algorithms have been adapted to work with sensor network's needs of low data rate and low power specification. These design modifications addresses mostly the functions of the MAC and the network layers and scale for topologies for less than few hundred nodes. The need for large environmental driven standards have not been specified, even though there has been public domain

software tools available such as TinyDB [9] which are still too complicated for an application developer due to its inherent design.

We design INSPIRE-DB framework from scratch to store the spatial and temporal characteristics of the ambient application parameters and shown that the pre-processing at the sensor level achieves fault-tolerance and the network data rate for periodic aggregation is as close to the theoretical limits of Slepian Wolf rate [5]. As the stored data need to be queried efficiently, the network topology is embedded into a labeled graph, which is mapped efficiently with data-centric storage. The aggregated data is routed efficiently using only the node labels of its neighbors without knowing the geographic node locations. The pre-processing used by INSPIRE-DB makes it possible to uniquely represent the smallest dimensionality of θ is the sparsity level [9] of the signal x under this model, making aggregation decision at the sensor level and not on routing protocol time periodicity. The paper is organized with an introduction to INSPIRE-DB framework in Sections II, III, and IV theoretical model for Distributed Source Coding and Compressed Sensing are compared for the fusion accuracy for same data-set. In Section V, DHT routing using data-centric storage is implemented. In Section VI, the simulation results are tabulated and scalability is discussed followed by Summary.

II. FRAMEWORK - INSPIRE DATABASE

Over main goal for an intelligent information processing, which is driven by *pre-processing* of the number of sensor measurements needed to be fault-tolerant and *post-processing* of the real-time data periodically to be *aggregated* and *routed* to its application label space for efficient *querying*.

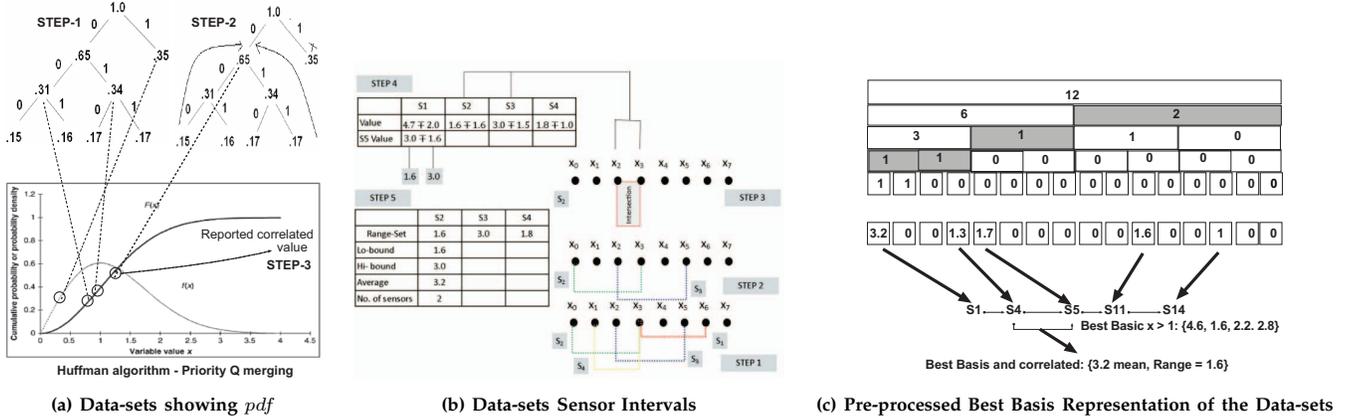


Figure 1. Sensor-centric Data Fusion During the Pre-processing Step for Different Representations of the Sensor Data-sets.

A. Theory Behind Sensor Processing

To design a model which allows to capture spacial and temporal nature of sensor data-sets and accurately represent the local measurement in time, we choose two data-centric fusion algorithms one, which is static and is based on the measured histogram [7] of the sensors and the other preserves the spatial and dynamic variations of the measured signal represented by using a dendrogram [7]. Both the algorithms needs atleast $(n - \tau)$ [2] sensor for reliably measuring the physical parameters without error. As the faulty nature of sensor networks can be described in number of sensors measurement n and the number of bits needed to representing the measurement b , is given in the equation (1). For a single correlated measured value the number of bits reduces to $b = 0$ for a large n due to high redundancy, if the same measured value is skewed as in the case of faulty sensors, then the bits needed to represent the measured values significantly increases shown in equation (2).

$$b_n \rightarrow 0 \text{ as } n \rightarrow \infty \quad (1)$$

$$b_n \cdot n \rightarrow \infty \text{ as } n \rightarrow \infty \quad (2)$$

1) *Histogram Model*: The measured values of the sensor X are independent identically distributed (i.i.d.) and can be represented as shown in equation.

$$f_x(x) = \sum_{i=-\infty}^{+\infty} f_x(i.i.d.) \quad (3)$$

Then the information content can be coded with a minimum number of bits is given by

$$H(S) = - \sum_{i=0}^i P(X_i) \log P(X_i) \quad (4)$$

The aggregation of the data-set is performed only

if the subset of the data-set values have a probability P_{max} of atleast ≥ 0.5 as shown in equation 5(a) and 5(b). This provides a way to determine locally if the sensors measurements are reliable as it represents the ground truth of the overlapping outputs of the sensors.

$$|P_{max}| = \begin{cases} \text{overlap}, & \text{for } P_{max} \geq \frac{n}{2} \\ \text{no-overlap}, & \text{for } P_{max} < \frac{n}{2} \end{cases} \quad (5a)$$

2) *Dendrogram Model*: A single measured signal of finite length, which can be represented in its sparse representation by transforming into its basis, then this technique is called the sparse basis in Compressed Sensing (CS) as shown in equation (6), of the measured signal, where Ψ_n^k is the best basis which is transformed from n to k and $(k \ll n)$. The technique of finding a representation with a small number of significant coefficients is often referred to as Sparse Coding. When sensing locally many techniques have been implemented such as the Nyquist rate [7], which defines the minimum number of measurements needed to faithfully reproduce the original signal. Using CS it is further possible to reduce the number of measurement for a set of sensors with correlated measurements.

$$x = \sum_{i=0}^n \vartheta(n) \Psi_n = \sum_{i=0}^k \vartheta(n_k) \Psi_{n_k}, \quad (6)$$

A sample data set which has been shown in Figure 1 (b) represents data from static sensors which have varying ranges. We use a dendrogram to represent the variable the range of the sensors S_1, S_2, S_3, S_4 in the interval $[0, 1, 2, 3, 4, 5, 6, 7]$.

The multi-sensor algorithm [6] forms the dendrogram tree is shown in Figure 1 (b) of correlated sensor intervals. In step1, the dendrogram is shown for the complete interval. In step2, the overlapping redundancies are removed by taking the best sensors intervals representing the overlapping intervals. Step3, finds

distributed geometrically with parameter a . The probability of successful transmission is given by probability mass function and lower-bound from Kraft inequality [1] $0 < a < 1$, to find a code minimizing:

$$\log \sum_{i \in \mathcal{X}} p(i) a^{l(i)} \quad (12)$$

$$P_{success} = a^{L_a(p,l)} \quad (13)$$

From the above equation (12), which are used by compression algorithm at the cluster heads[1] to assign the prefix code optimally, the quantitative sensor data optimization is due to coefficient a , the periodic distribution of i.i.d. values over transmit time in P , the probability of success.

IV. DISTRIBUTED COMPRESSED SENSING (DCS)

DCS allows to enable distributed coding algorithms to exploit both intra-and inter-signal correlation structures. In a sensor network deployment, a number of sensors measure signals that are each individually sparse in the some basis [7] and also correlated [2] from sensor to sensor. If the separate sparse basis are projected onto the scaling and wavelet [7] functions of the correlated sensors(common coefficients), then all the information is already stored to individually recover each of the signal at the joint decoder. This does not require any pre-initialization between sensors.

A. Joint Sparsity Representation

For a given ensemble X , we let $P_F(X) \subseteq P$ denote the set of feasible location matrices $P \in P$ for which a factorization $X = P\Theta$ exists. We define the joint sparsity [9] levels of the signal ensemble as follows. The joint sparsity level D of the signal ensemble X is the number of columns of the smallest matrix $P \in P$. In these models each signal x_j is generated as a combination of two components: (i) a common component z_C , which is present in all signals, and (ii) an innovation component z_j , which is unique to each signal. These combine additively, giving

$$x_j = z_C + z_j, j \in \forall \quad (14)$$

$$X = P\Theta \quad (15)$$

B. Distributed Fused Parameter Dictionary

The sample sensor measurements of Figure 1 (b) is transformed in Figure 1 (c) to show all its possible basis representations. The cost-function [7] searches to find an optimal (grey rectangles) best basis matching the least number of coefficients to represent the signal without overlaps. The lowest range is calculated by selecting consecutive significant coefficients (1.3, 1.7), which determines the maximal overlap for the sensor

intervals. This best basis dictionary is stored in the hashed location of the application's search tree.

V. DATA-CENTRIC ROUTING

DHT was inspired by the design of the P2P internet in CHORD [4] and Tapestry [4]. Post-processing of the fault-tolerant data-sets are aggregated using data-centric routing. To distribute the data to its corresponding key space, DHT based mapping is used. In Sensor networks, multi-hop and best effort routing, the protocols select routes from many data sources to a central base station. These categories of routing depends on geographic or node's location information, in the case of data-centric routing the fused sensor data is mapped to a virtual key space making routing decisions distributed compared to short-path routing.

A. Fused Objects

The distributed sensor processing uses a notion of local and distributed dictionary of the parameters it has currently sensed and for the new fused values. The local dictionary consists of all the occurring atoms of the individual measurements as shown in equation (7) and its distributed version represents the maximal overlap of the multi-sensor ranges as shown in equation (15). On completion of the pre-processing step, the measured parameter will have fused value and it's minimum and maximum ranges. The global dictionary is stored in a compressed form and can be used to recreate the individual measured signals without loss. This also allows for quick lookup when its application defined parameters is queried, it returns the stored value seen and its range.

B. Aggregation Tree

A root node is selected from the node topology, which is conveniently placed for querying and a minimal spanning tree is formed using the neighborhood table as shown in Figure 2 (b). The tree nodes are given a key value using a uniform hash function over the total number of nodes. A simple tree can be formed using a flooding protocol [1] with the nodes wireless range and assigning nodes the hop level starting from the root and so on to form a tree with L-levels, as shown in Figure 2 (c). These levels are labeled to store the aggregated data-sets.

C. Key Space

The Sensor network application is measuring the parameters $X_1, X_2, X_3, \dots, X_m$, then the fused physical measured values at the sensors are aggregated at a rate $\Upsilon_1, \Upsilon_2, \Upsilon_3, \dots, \Upsilon_m$. The data-centric key map of X is given $f(X)$ and stored in the aggregation tree corresponding to node's key. To query the root node for the

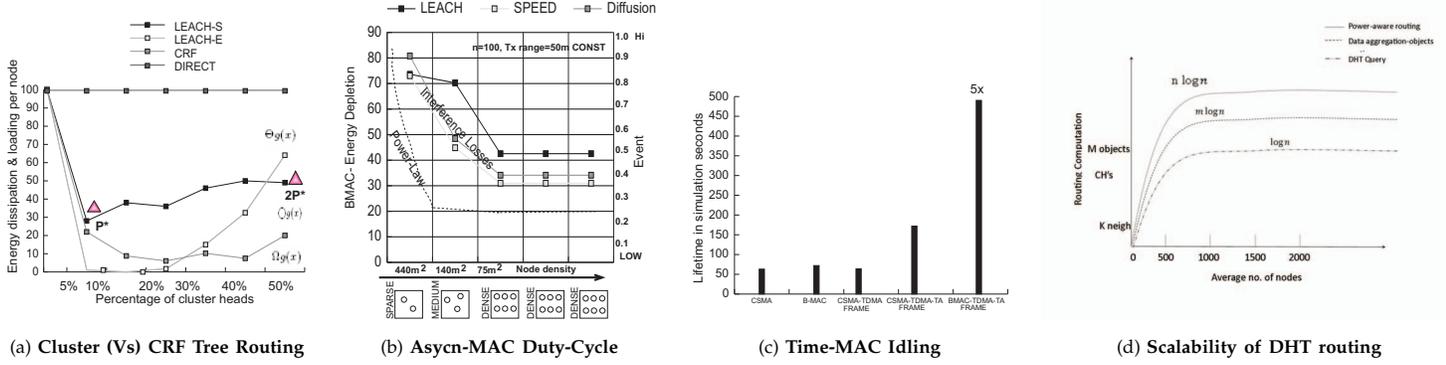


Figure 3. Topology Control Protocols for Short-Path, Power-aware and DHT Routing for Large Sensor Networks.

MAC LESS Cross Layer Simulation Bayesian fault rate Fixed Energy Model	$\omega_1 = \omega_2$	$\omega_1 \neq \omega_2$	Assumptions	MAC LOSSES Protocol Simulations Bayesian fault rate Renewable Energy Model	CSMA	BMAC	Assumptions
LEACH Fixed Energy Node failures	27% $\omega \leq 20\%$ *	41% $P = 2P^*$ $P \leq 2P^*$	Residual Energy Optimal config BE Error	Clustering Best Effort Renewable E-Model Channel Error Model	50% better 22% better * $P = P^*$	55% better 35% better $P \geq 2P^*$ $P \leq 2P^*$	No-MAC No-MAC Node failures. Link errors.
Network Scalability Model Pre-Processing Post-Processing	$O(c \lg n)$ $O(\lg n)$	$O(n^2 d^2)$ $O(n^2 d^2)$	Sensor Fusion Data Aggregation	Cross-layer Protocols Pre-Processing Post-Processing	$P < P^*$ $P = P^*$	$P = P^*$ $P \leq P^*$	INSPIRE-DB Link Quality Multi-hop protocols

Table I

SENSESIM: SIMULATION TEST-BED FOR POWER-AWARE LIFETIME MODELS

value of X , it needs to compute $f(X)$ and get the key and send the message to the same node corresponding to that key. To store a data value its hash corresponds to a given tree level (parent node) and a matching offset (child node), which is shown as an embedded tree in Figure 2 (c).

D. Prefix-Routing

The routing from $\langle \text{SRC}, \text{DST} \rangle$ overlay messages to the destination ID digit by digit as shown in Figure 2 (c) (e.g., $6*** \Rightarrow 62** \Rightarrow 62A* \Rightarrow 62AD$, where $*$'s represent wild-cards). This addressing scheme is similar to longest prefix CIDR specification. As the queryable address space is evenly distributed for n nodes, the data-sets which are replicated are also evenly distributed to any specific DHT source entry. The complexity of the DHT routing can be computed as $O(\lg n)$ hops from any node to a given label location in the tree. Due to dynamic changes in topology if some nodes are not available the next node in the same level is assigned the closest $\langle \text{KEY}, \text{VALUE} \rangle$ pair thus, effecting locally the nodes which belong to the same aggregation bin, making it resilient to failures.

E. Node-to-Node Routing

From the above Pre-fix routing a node can communicate and get its data by knowing the labels of its neighbors, similarly if a node has a static identifier n , if a message needs to be sent to the node, then its label $L(n)$ can store its label using $f(n)$ the same way it stored the data. If another node needs to communicate with n then it looks up the label using $f(n)$ in the DHT and retrieves $L(n)$. This is used to route messages between the nodes in the network topology.

F. Low Embedding Overhead

For DHT, to locate the stored value for a given query key the data need to be also copied to the application's mapped KEY label according to the routing tree information. We need to embed the aggregation tree into the nodes topology in a distributed manner, without incurring added data aggregating overheads. Routing once the pre-processed fused levels are available it is hashed, the storage locations are determined by the hash value in the DHT address space and if a live node exists for the same KEY index, then it becomes the root for all the aggregated data set, otherwise one of the children stores the data value. As DHT hops may not

translate into reachable nodes in the deployed sensor network topology, due to this the routing table needs to maintain a path when the edges of the graph are not reachable. This embedding [10] of the tree levels is called dilation [10]. The routing table maintains j levels of upstream and downstream pointers, if c is the number of neighbors and IDs are generated of base β then the total number of pointers is $c\beta$ which differ by at least 3 hex digits.

VI. SIMULATION

We categorize the sensor routing algorithms into central routing (LEACH) [1], Geographic routing (SPEED) [1], Best effort routing (B-MAC) [8], Data dissemination (Directed Diffusion) [8], Tree based (CRF) [1] routing and Tree based routing with sensor pre-processing and aggregation [INSPIRE-DB]. The simulation [3] is based on cross-layer energy performance for networks and MAC layers for each of these categories of protocols. As the lifetime is an important performance measure, we simulate for lossless medium and for MAC with radio propagation losses (collision, Idle and Sleep) [8]. The testbed results are categorized for MAC and MAC-Less simulations. Power-aware routing can be achieved by actively load balancing the nodes and scheduling the MAC to save energy during idle as shown in Figures 3 (a), (b) and (c). None of the non-tree protocols implement label based DHT routing and relies on node's geographic information and does not optimize on features like data-centric storage. Due to having node storage and sensor fusion fault-tolerant pre-processing using INSPIRE-DB only fused data need to be hashed to its corresponding key location address and stored, minimizing network wide traffic.

VII. SCALABILITY TO LARGE NETWORKS

Next we examine how well INSPIRE-DB scales with the size of the network with constant number of sensed parameters. As the framework needs to fuse the physical values from the sensor outputs and then aggregate the data to its data-centric locations in a distributed way using equations (8) and (15), which enables applications to look-up the KEY for the same node in the DHT and send a query message. The computation of routing complexity for sensor-centric as well as data-centric are shown in Figure 3 (d) for networks from 10 to 2000 nodes. As DHT uses a tree based routing scheme, the search overhead grows logarithmically with the size [10] of the network. The average query latency to locate the data in the stored nodes for a large network remains constant as the query path for node-to-node increases.

VIII. SUMMARY

The performance improvements for each cross-layer and its application needs are compared from the results of the testbed with 100 to 5,000 nodes as shown in Table (I). The INSPIRE-DB application framework uses pre-processing making it sensor centric. The framework is fault-tolerant due to the pre-processing criteria of $(n - \tau)$ and data-centric aggregation minimizes the amount of data, which is sent to the central base station by only returning data-sets for the query range requested.

REFERENCES

- [1] Vasanth Iyer, S. Sitharama Iyengar, Garimella Rama Murthy, N. Balakrishnan, and Vir Phoha. Distributed source coding for sensor data model. In Proc. Third International Conference on Sensor Technologies and Applications SENSORCOMM, 17-21 June. 2009.
- [2] Bhaskar Krishnamachari, Member, IEEE, and S. Sitharama Iyengar, Fellow, IEEE. Distributed Bayesian Algorithms for Fault-Tolerant Event Region Detection in Wireless Sensor Networks, IEEE Transactions on Computers, Vol. 53, No. 3, March 2004.
- [3] Vasanth Iyer, S. Sitharama Iyengar, Garimella Ram-murthy, and Mandalika B. Srinivas. SenseSIM: Sensor Network Simulator, ISSNIP, Melbourne, Australia, 2009.
- [4] Krishna Gummadi, Ramakrishna Gummadi, Steve Gribble, Sylvia Ratnasamy, Scott Schenker, and Ion Stoica. The impact of DHT routing geometry on resilience and proximity. In Proc. of SIGCOMM (Karlsruhe, Germany), ACM, pp. 381394, Sep 2003.
- [5] Slepian, D. Wolf and J. Noiseless coding of correlated information sources, 1973.
- [6] Richard R. Brook, and S. S. Sitharama Iyengar. Robust Distributed Computing and Sensing Algorithm, ACM, 1996.
- [7] Arne Jensen and Anders la Cour-Harbo. Ripples in Mathematics, Springer Verlag, 2001. 246 pp. Softcover ISBN 3-540-41662-5.
- [8] Joseph Polastre, Jason Hill, and David Culler. Versatile low power media access for wireless sensor networks SenSys: Proceedings of the 2nd international conference on Embedded networked sensor systems, ACM, 95-107, 2004.
- [9] Dror Baron, Marco F. Duarte, Michael B. Wakin, Shriram Sarvotham, and Richard G. Baraniuk. Distributed Compressive Sensing. In Proc: *Pre-print*, Rice University, Texas, USA, 2005.
- [10] James Newsome and Dawn Song. GEM: Graph Embedding for Routing and DataCentric Storage in Sensor Networks Without Geographic Information. Proceedings of the First ACM Conference on Embedded Networked Sensor Systems (SenSys). November 5-7, Redwood, CA, 2003.