

## Distributed Source Coding for Sensor Data Model

Vasanth Iyer  
IIIT Hyderabad  
A.P., India 500032  
vasanth@research.iiit.ac.in

Sundaraja Sitharama Iyengar  
Louisiana State University  
Baton Rouge, USA 7081  
iyengar@csc.lsu.edu

Rama Murthy Garimalla  
IIIT Hyderabad  
A.P., India 500032  
rammurthy@iiit.ac.in

Srinivas Mandalika  
IIIT Hyderabad  
A.P., India 500032  
srinivas@iiit.ac.in

**Abstract**—We measure reliability in sensor networks which are dependent on limited resources of individual sensor nodes such as battery capacity, transmission range and channel interference due to simultaneous wireless transmissions. From the initial simulation it is estimated that the routing errors using a distributed algorithm for a large network is less susceptible to failures when compared to using a table driven routing algorithm. To further address other influencing factors which are not related to resource allocation or routing of the sensor network we study the correlated issues, which makes sensor network unique to the categories of wireless network applications. The simulation results show that due to 1-bit-mask accuracy and the CDF codes used to represent measured values in the decoder buffer is fault-tolerant and also increases the communication rate by 70% due to information redundancy within a sensor cluster.

**Keywords**—Sensor Data Reliability, Slepian & Wolf Coding, Cosets, Huffman Trees, BER, Bayesian Error.

### I. INTRODUCTION

Sensor networks are deployed in a dense configuration due to its limited radio range and fixed non renewable energy resources due to computational/networking characteristics of sensor networks. To collaboratively use the limited resources distributed algorithms, select a single node which transmits serially using its UART pre-processed sensed data information using many local resources. As the cost of radio transmission is much more than local sensing, the sensor network uses two different topologies to address the energy cost at the cross-layer stack. The network layers uses the upper layers assuming MAC layer abstraction to optimally pick cluster heads by using a fixed probability density function (pdf) of a network resource at the node, such as, remaining battery energy. This type of pdf is power-aware as it uses a collaborative function to minimize over use of network resources thus avoiding pre-mature node failures.

The MAC layer uses a k-neighborhood distance algorithm to find other nodes within its own limited range and uses a multi-hop schedule to the specific data transmitting node. This scheduling allows multi-hop nodes to use sleep cycles and lower their energy consumption while idling. These multi-hop algorithms use low-power listening and use a preamble to wake up nodes, sleep cycles when the transmitter is completely off and traffic based preamble to synchronize nodes to receive the data payload.

If  $\theta_1$ ,  $\theta_2$ ,  $\theta_3$  are the data values of a parameter such as residual energy, observed values by the sensors, as large scale sensor deployment are a dense deployment as the reading are correlated only an average  $\theta_1$  needs to be transmitted. As the clustering is based on the network layer which optimizes on radio range and not the sensing region it always is approximated and corrected using some training samples using less number of bits to be transmitted, this is the fundamental design based on power-aware data model.

In the MAC layer which polls the channel to check for any activity while receiving and during transmitting to avoid collision and uses best effort QoS for the messages to be forwarded. The data sensing nodes are single hop, while forwarding nodes are multi-hop nodes.

The data values which are forwarded are discrete and updated according to some trend in the data. Some measured values may be changing more quickly than others creating different traffic patterns that are data driven. The multi-hop nodes do not have any sensors and act like routers which uses best effort QoS and constantly adapts its polling depending on the data trend, this is fundamental to the design of polling the channel, which uses on-demand traffic predictions. Model implementation assumes  $\theta_1$ ,  $\theta_2$ ,  $\theta_3$  are always transmitted when changes happen and typically it is re-transmitted at a constant rate of 10 minute intervals keeping the channel polling of a set of nodes to guarantee the QoS. Figures 1(a,b) illustrates the Bayesian classifier for pdf based clustering and multi-hop based passive clustering. For the theoretical and mathematical proofs please refer to chapter 2,4 in the mentioned reference [7].

This paper builds from previous work [1,2] and extends the two dimensional Bayesian model [7] to optimize on power-aware routing algorithms in representing sensor network. The routing algorithms are implemented at the network layer which have known density of nodes by using prior selection and at MAC layer which have unknown node densities due to limited transmission range. The Bayesian classifiers [7] which are specific to the routing topology uses features to maximize the lifetime of the sensor network and minimize on sensor faults. This Bayesian classifier helps in predicting the theoretical fault rate bounds by knowing the node densities validated also through real simulation.

In section II, the sensor data model is described with respect to sampling and compression needs at the cluster heads. In section III, the Source coding rate is introduced for correlated sources using error corrected codes. In section IV, the scalability of the sensor network is modeled using Power Law and Bayesian Classifier and how it effects distributed clustering and passive clustering routing. In Section V, a distributed algorithm is simulated without MAC to find error bounds for a large-scale deployment. Section VI, uses battery model at the physical layer, which uses the cross-layer energy

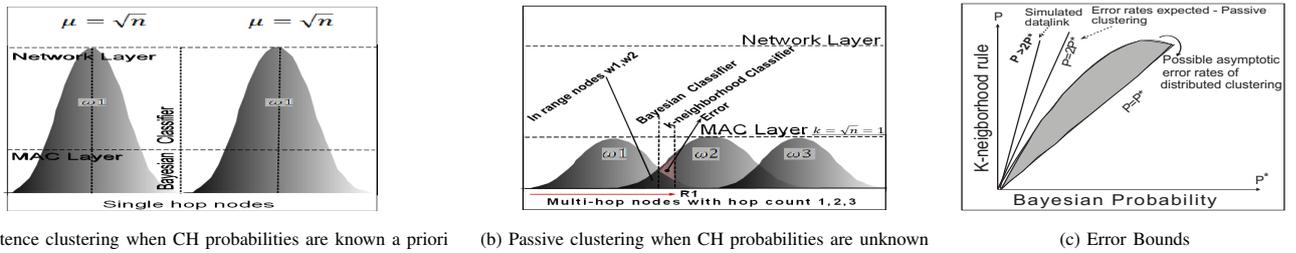


Fig. 1. Estimation of CH selection error and MAC layer routing using Bayesian distributed rule

model with a standard simulator using crossbow mote energy model to analyze lifetime for various routing. Section VII uses fault recognition pre-processing using Huffman probability trees and measures the bit rate errors in the networks. Section VIII uses probabilistic algorithms to calculate a local variance, which helps in determining the transmission success rate.

## II. DATA MODELS

### A. Probability model

Where  $d$  is the distance to transmit between sensors  $i$  to sensor  $j$ . We designed the compression algorithm for a large distributed sensor network with a desired channel rate, a fixed length code to represent real sensed values at the encoder using  $(c,k,d)$ , Where  $c$  is the code,  $k$  is the length and  $d$  is the distance from the average  $P_{Max}$  cluster heads. This technique which allow using less number of bits to represent the newly encoded data is sent to the decoder by sharing the expected local value at both ends.

As a rule, compression algorithms use a probability model based on the entropy of the source. Iyengar [3] defined a Bayesian fault-tolerant algorithm in sensor network using an abstract sensor which can be *tamely faulty* and *widely faulty*. For larger sensor network deployment, this model helps predict the error bounds in terms of the varying sensing values. In this paper we adapt the Bayesian rule [7] to select cluster heads for known node density and extend it to find the upper bounds related to unknown densities for solving the optimal routing problem at the network layer in sensor networks. The latter is more relevant for renewable energy resources [4]. Where the number of active sensors connected to the network is not known, at any given time.

Entropy of general sensing source is a sequence of length  $n$  from the source and is given by

$$H(S) = \lim_{n \rightarrow \infty} \frac{1}{n} G_n, \text{ where} \quad (2.1)$$

$$G_n = - \sum \sum \dots \sum P(X_1 = i_1, X_2 = i_2 \dots X_n = i_n) \log P(X_1 = i_1, X_2 = i_2 \dots X_n = i_n)$$

In sensors where each element in the sequence is independent and identically distributed (i.i.d.), with this statistical model, we can modify the entropy of the first order to equation (2.1)

$$H(S) = - \sum P(X_1) \log P(X_1) \quad (2.2)$$

### B. Aggregation model

If the cluster size in  $n$ , given this density of clustering, the entropy of data aggregation [8,9] from equation 2.2. In a lossless mode if there are no faults in the sensor network, we can show that the highest probability given by  $P_{Max}$  is ambiguous if its frequency is  $\leq \frac{n}{2}$  otherwise it can be determined by a local function.

### C. Local $P_{max}$ functions

Provides a way to determine the local filter value from the probability distribution used by compression algorithms.

$$|P_{max}| = \begin{cases} \text{local}, & \text{for } P_{max} \geq \frac{n}{2} \\ \text{global}, & \text{for } P_{max} < \frac{n}{2} \end{cases} \quad (2.3a)$$

### D. Slepian & Wolf theorem

The Slepian-Wolf rate [11] region for two arbitrarily correlated sources  $x$  and  $y$  is bounded by the following inequalities, this theorem can be adapted using equation (2.2)

$$R_x \geq H\left(\frac{x}{y}\right), R_y \geq H\left(\frac{y}{x}\right) \text{ and } R_x + R_y \geq H(x, y) \quad (2.4)$$

If the correlated sources are differing by a few bits, the possible number of codewords can be represented as  $2^m$  where  $m$  = no. faulty bits [10]. In our case  $m=2$  as the parameters are distributed whilst collected locally at the cluster head.

## III. COMPRESSION RATE

### A. Distributed source coding with side information

In sensors networks several measured values are sensed in a distributed manner and these are aggregated according to the users query. The goal of all the encoder is analogous to the previous section where it uses cosets. Equations (3.1-3.4) illustrates the bin formation to reduce the overall bits needed for transmission. Considering the case of distributed sensing application, the encoder is further designed with a machine learnable redundancy range which is specific to each and every application. This mutually redundant measured range is correlated with sensors which are in the same wireless range and connected to a parent. This information, also called side information is shared with the decoder. Owing to side information, even lesser number of bits are needed to represent the changing values coming from each cluster heads transmitting to the joint decoder. Encoder and decoder have access to the side information  $Y$ . which is correlated to  $X$

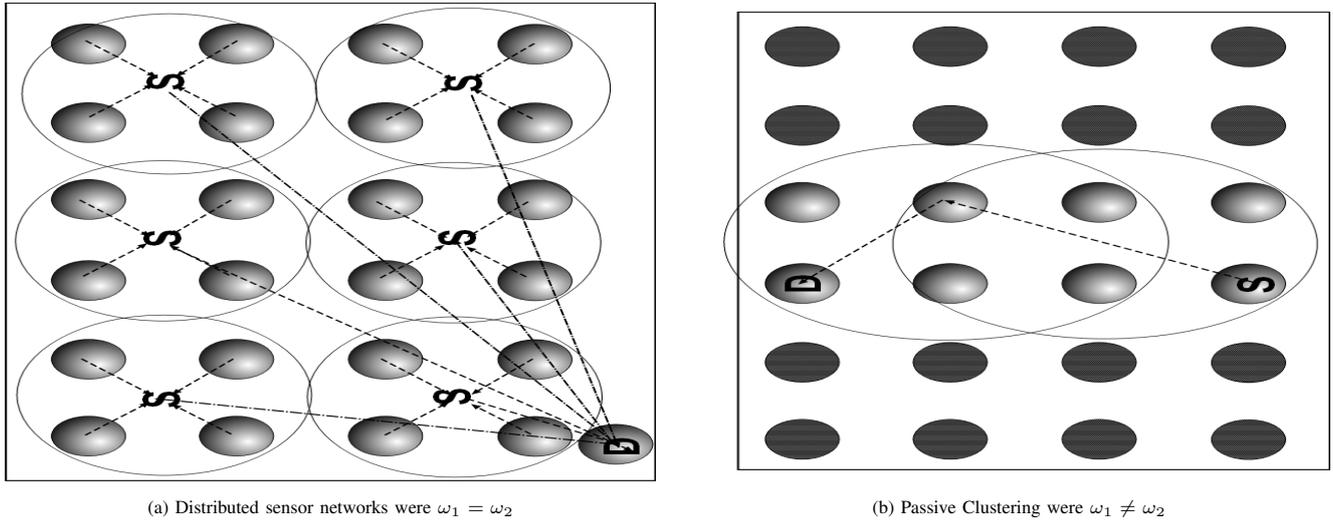


Fig. 2. (a) Distributed sensor networks with LEACH single-hop nodes (b) Passive clustering algorithms with multi-hop nodes

and can be represented by the equation 2.3(a,b). According to the Slepian-Wolf Theorem [11], established in 1971, that the number of bits needed by using the theorem is lesser, as shown in figure 4(b), than the total entropy for both the two arbitrarily correlated sources  $H(x)$ ,  $H(y)$ .

$$\begin{vmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \end{vmatrix} = 00 \quad (3.1)$$

$$\begin{vmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \end{vmatrix} = 10 \quad (3.2)$$

$$\begin{vmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \end{vmatrix} = 01 \quad (3.3)$$

$$\begin{vmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \end{vmatrix} = 11 \quad (3.4)$$

#### IV. FAULT RATE

Large deployment of sensor network that use an efficient distributed algorithm to select cluster heads which allows to extend the lifetime [5] to function without faults. The fault rate of such an algorithm can be defined as the residual percentage of good sensor when the network incurs faults due to resource drain. This is typically referred to as the sensor networks residual energy, if the fault rate is higher the cluster head selection algorithm is less optimal. The two dimensional simulation model is expressed in figure 2 for distributed and passive cluster based routing. In the paper the fault rate is measured for both the cases for algorithm complexity, multi-hop dependency, MAC layer losses and Bit error rates.

##### A. Estimate of the sensed value for known densities

**Theorem IV.1** A power law is any polynomial relationship that exhibits the property of scale invariance. The most common power laws relate two variables and have the form.  $PowerLaw = f(x) = ax^2 + o(x)^2$

*Proof:* The function  $f(x)$  is represented as function of transmission distance from the cluster heads to a sink location,

$f(d)$ , where  $d$  is the distance to transmit between sensors  $i$  to a multihop sensor  $j$  towards the sink in increasing distance, from this we get the Power rule [5] based on the distance  $d$  of nearest sensor to the farthest away sensor, substituting in the above theorem IV.1 and summing up the total energy required for all transmissions within one meter, two meters, three meters, four meters and extending up to  $(d - 1)$  meters to a progressive sequence in equation (4.1).

$$PowerLaw = 1^2 + 2^2 + 3^2 + 4^2 + \dots + (d - 1)^2 + d^2 \quad (4.1)$$

To sum up the total energy consumption we can write it in the form of Power Law equation (4.1.1)

$$PowerLaw = f(x) = ax^2 + o(x)^2 \quad (4.1.1)$$

Substituting  $d$ -distance for  $x$  and  $k$  number of bits transmitted, we equate as in equation (4.1.1).

$$PowerLaw = f(d) = kd^2 + o(d)^2 \quad (4.1.2)$$

$o(d)^2$  is an asymptotically small function of  $d$ , Taking Log both sides of equation (4.1.2),

$$\log(f(d)) = 2 \log d + \log k \quad (4.1.3)$$

Notice that the expression in equation (4.1.2) has the form of a linear relationship with slope  $k$ , and scaling the argument induces a linear shift of the function, and leaves both the form and slope  $k$  unchanged. Plotting to the log scale. ■

**Corollary IV.2** Properties of power laws - Scale invariance: The main property of power laws that makes them interesting is their scale invariance. Given a relation  $f(x) = ax^k$  or, indeed any homogeneous polynomial, scaling the argument  $x$  by a constant factor causes only a proportionate scaling of the function itself. From the equation (4.2.1) we can infer that the property is scale invariant even with clustering  $c$  nodes in a given radius  $k$ .

<b>Randomized CH Selection Scheme</b>
<b>Generate a random number</b> $x \in (0, \%CHs)$ <b>Calculate</b> $g_i(x) = P(\omega_i \mathbf{x}) = \frac{p(\mathbf{x} \omega_i)P(\omega_i)}{\sum p(\mathbf{x} \omega_j)P(\omega_j)}$ <i>if</i> $x = rand(x)$ <b>if</b> $x \leq \%$ <i>then</i> $CH_i = x$ , <i>else</i> $CH_i = false$
<b>Threshold CH Selection Scheme</b>
<b>Obtain the sensors residual energy</b> $S_j$ <b>for all</b> $N_i$ <b>neighbors of node i</b> <b>Calculate if</b> $\theta \leq S_j$ $g_i(x) = P(\omega_i \mathbf{x}) = \frac{p(\mathbf{x} \omega_i)P(\omega_i)}{\sum p(\mathbf{x} \omega_j)P(\omega_j)}$ <i>if</i> $x \leq \theta$ <i>if</i> $x \geq \theta$ <b>then</b> $CH_i = x$ , <i>else</i> $CH_i = false$
<b>Optimal Zone based caching Scheme</b>
<b>Divide the sensors into three zones</b> <b>Use the middle zone as CHs caching</b> <b>Calculate</b> $g_i(x) = P(\omega_i \mathbf{x}) = \frac{p(\mathbf{x} \omega_i)P(\omega_i)}{\sum p(\mathbf{x} \omega_j)P(\omega_j)}$ <i>if</i> $x = Constant$ <b>Use optimal settings from the above two cases for % of CHs with use count</b>

(a)

Fig. 3. Cluster head selection for power-aware routing in large sensor networks

*Proof:*

$$f(d) = kd^2 + o(d^2) \quad (4.2)$$

$$f(cd) = k(cd^2) = c^k f(d) \alpha f(d) \quad (4.2.1)$$

From the equation (4.2.1) we can infer that the property is scale invariant even with clustering  $c$  nodes in a given radius  $k$ . This is validated from the simulation results [6] obtained in Fig 4 (a) which show optimal results (minimum loading per node [6]) when clustering is  $\leq 20\%$  as expected in theorem 1. ■

**Theorem IV.3** Theorem 3. *CH Error Rate - Local:* *If two classes have the same covariance, where*  $p(x|\omega_j) \approx N(\mu, \Sigma)$ ,  $j = 1, 2$ .

If prior probabilities are equal and the event are independent across clusters, the Bayes model minimizes according to the input distribution and the error rate is given by

$$P(e) = \frac{1}{\sqrt{2\pi}} \int_{r/2}^{\infty} e^{-u^2/2} du \quad (4.3.1)$$

*Proof:* The simulated algorithms such LEACH use the knowledge that the nodes which are sensing are correlated and have known densities such as cluster size and radio range. The sensed values are i.i.d distributed and their variance  $\neq 0$ . The underlying model can use the error rate for a cluster as  $\frac{1}{N}$  and estimated value  $\theta$  which is random value of the nodes residual power in this model and is defined by

$$r^2 = \int_{i=1}^d \left( \frac{\mu - \mu}{\sigma_i} \right)^2 \quad (4.3.2)$$

where  $r^2$  is the radio range between nodes calculated by using Mahalanobis distance [7].

$$P = P^* \quad (4.3.3)$$

Which is the classifiers error bounds and is the lower bound Bayesian estimate. ■

**Theorem IV.4** Theorem 4. *Multi-hop Error rate - Global:* *When*  $P(\omega_m|x)$  *is close to unity, the nearest-neighbor selection is almost always the same as Bayes selection. This is, when the minimum probability of error is small close to*  $1/c$ , *so that all classes are essentially equally likely, the selection made by the nearest-neighbor rule and the Bayes rule are rarely the same, but the probability of error is approximately*  $1 - 1/c$  *for both and is bounded by.*

$$P \leq 2P^* \quad (4.4)$$

*Proof:* We recall that the Bayes decision rule minimizes  $P(e)$  by minimizing  $P(e|x)$  for every  $x$ . If  $P^*(e|x)$  be the minimum possible value of  $P(e|x)$ , and  $P^*$  be the minimum possible value of  $P(e)$ , as the probability of an event is conditional across neighbors is can be represented as  $P^2$ . The error probability is given by

$$P^*(e) = 1 - \sum_{i=1} P^2(\omega_i|x) \quad (4.4.1)$$

and from the previous theorem convergence of the Nearest Neighbor( $k$ ), then the error rate is  $\frac{1}{k}$ . The reusable probability at the CH for passive(local) clustering compared to distributed clustering is  $\frac{1}{k} > \frac{1}{N} \leq 1$

$$P^*(e) \approx 2(1 - P(\omega_i|x)) \quad (4.4.2)$$

$$P \leq 2P^* \quad (4.4.3)$$

Equation 4.3.3 and 4.4.3 differentiates the error rate in terms of residual energy index(rei). ■

## V. ANALYSIS OF FAULT RATE CH-ALGORITHMS

### A. Estimate of the sensed value for known densities

The simulated routing algorithms such LEACH-S [8], LEACH-E [6] and CRF [6] as described in the above table use the knowledge that the nodes which are sensing are correlated and have known densities such as cluster size and radio range.

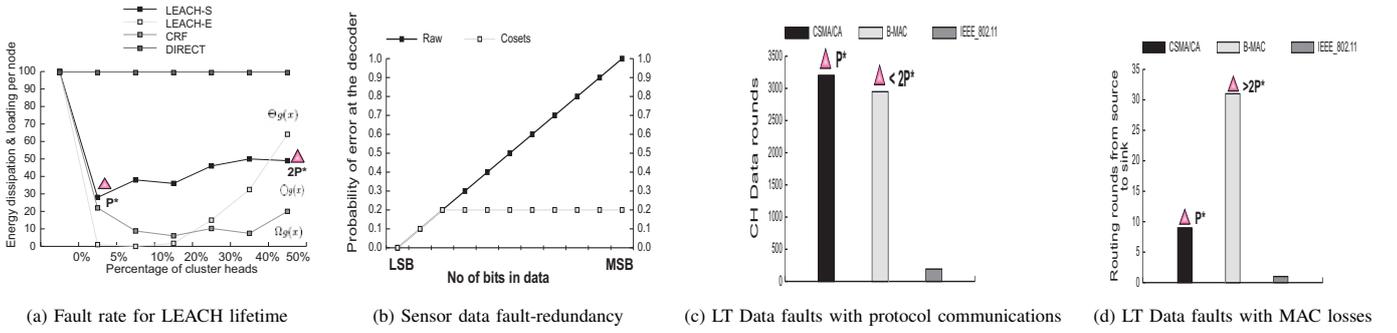


Fig. 4. Simulation results of fault analysis of WSN routing and data aggregation algorithms

The sensed values are i.i.d distributed and their variance  $\neq 0$ .

The underlying model uses different ways to select the cluster heads to minimize the error rate. When the sensor faults happen due to fixed energy resources at the cluster head the total energy unused at the end of its lifetime is the residual rate[2], the routing algorithms tries to minimize this error criteria. As this model uses the network layer as shown in figure (2a) and the only dependant variable is the fixed lifetime model [5].

The complexity of the algorithm can be defined by using the standard implementation of the LEACH distributed algorithm and its power-aware variations, see table (figure 3).

$$\bigcirc g(x) = f(x) : 0 \leq f(x) \leq cg(x) \quad (5.1)$$

$$\bigtriangleup g(x) = f(x) : 0 \leq cg(x) \leq f(x) \quad (5.2)$$

$$\bigodot g(x) = f(x) : 0 \leq c_1g(x) \leq f(x) \leq c_2g(x) \quad (5.3)$$

Complexity of the routing algorithms for LEACH is shown in equation (5.1), LEACH-E equation (5.2) and CRF equation (5.3). In the next section we will use only the lower layer such as power-ware MAC and estimate the multi-hop routing errors. In this case the model is not dependant on the fixed energy resources and only dependant on k-neighborhood rule it uses to find its multi-hop nodes as shown in figure (2b).

As the node probability are not known a priori the error rates are much higher than the persistence clustering.

## VI. SIMULATION

### A. Results from the network layer

Simulation models large number of nodes and calculates the lifetime when sensor faults are more likely to happen, the table shows(see Table I) number of cluster heads and the fault rate for distributed clustering and passive clustering in figure 4. Simulation results confirms the fault rate is network size invariant and converges to the optimal values derived in theorem 1 and 2.

### B. Results from the MAC layer

When node densities are not know in advance due to node failures or unscheduled polling and other characteristics of sensor due to its dependence in fixed resources. The problem due to this is for data transmitting nodes needs to find a near

TABLE I  
SUMMARY OF NOTATIONS FOR ANALYSIS OF ROUTING FAULT-RATE

Symbols	Definition
$N$	Total number of deployed nodes
$n$	Number of nodes in the cluster
$\mu$	Density of the class
$P_{MAX}$	Bayesian class rule
$R_x, R_y$	Entropy of correlated sources
$R, r$	Radio Range
$P$	K-neighborhood fault probability
$P^*$	Bayesian probability
$\omega$	Bayesian classes
$S$	Data source node
$D$	Destination node
$\theta$	Nodes residual energy
$CH$	Cluster head
$P(\omega_i    \mathbf{x})$	Conditional probability
$P(\mathbf{x}    \omega_i)$	Class conditional probability

neighbor in a deterministic way by which it can build a passive cluster to multi-hop its data. This uses minimal clustering overhead as it does not use the upper layers during communication synchronization. The behavior of the k-Nearest-Neighbor rule [7] will be directed by in our simulation a two-dimensional node distribution of  $n \geq 100$  where node density has one or less neighbors. The unconditional average probability of error occurring will be found over all nodes positioned at coordinates specified by x:

$$P^*(e) = \int P(e|x)p(x)dx \quad (6.1)$$

The convergence of the nearest neighbor for distributed clustering and passive clustering are derived, the distributed clustering case is

$$P = P^* \quad (6.2)$$

For passive clustering is given by

$$P = 2P^* \quad (6.3)$$

As shown in figure 2(a) where lower bound for LEACH-S when it becomes faulty and the remaining residual energy using the cross-layer simulator is  $P(e)=0.27\%$  which is the fault rate.

### C. Results from the MAC layer using a propagation model

In the case of passive clustering when node density  $p=0.1$  or using the k-neighborhood rule as shown in 2(b) the node densities are unknown in this case due to high likely-hood of

MAC LOSS 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	$\omega_1 = \omega_2$	$\omega_1 \neq \omega_2$	Assumptions
	LEACH Fixed Energy Node failures(renewable lifetime)	0.27% $\omega < 20\%$ x			0.41% $P = 2P^*$ $P \leq 2P^*$	$2^m$ where m=2 faulty bits Optimal config, Theorem 1,2 Errors due to unbalanced nodes	

Fig. 5. Simulation test-bed for power-aware lifetime models

faults. The protocol simulation results are show in table (figure 5) that the upper bound has error rate of  $P(e)=0.41\%$  which converges to the proof derived in theorem 3 and theorem 4 and the upper bound in figure 1(c), chapter 4 on non parametric techniques [7].

In the previous case MAC abstraction is used which does not take into account the propagation losses and protocol retries at the MAC level. To simulate the wireless channel we use GlomoSIM [4] bit error rate(BER) simulator and implement the routing algorithms for multi-hop cases. The routing algorithm implemented is SPEED which as shown in table (figure 5) is a geographic routing algorithm which uses two dimensional coordinate space to calculate the path from the node coordinates.

Many runs into the protocol simulation suggest that the radio characterization for CSMA [4] and B-MAC are comparable, figure 4(c) when the node densities are known.

The radio characterization for CSMA [4] is prone to faults when compared to B-MAC, figure 4(d) when using in multi-hop modes where the node densities are unknown. The protocol performance results show that the data packets received during useful lifetime is 3X times better in B-MAC when compared to CSMA and error rates are  $P \geq 2P^*$  higher than the theoretical Bayesian limit [7] of  $P = 2P^*$  as derived in theorem 3 and theorem 4.

## VII. FAULT RECOGNITION

Without loss of generality, we will assume a bitwise mask model of the data in which a particularly large value is considered unusual(MSB), while the normal reading is typically a low value. If we allow for faulty sensors, sometimes such an unusual reading could be the result of a sensor fault, rather than an indication of the event. We assume environments in which event readings are typically spread out geographically over multiple contiguous sensors. In such a scenario, we can disambiguate faults from events by examining the correlation in the reading of nearby sensors.

Let the real situation at the sensor node be modeled by a binary variable  $T_i$ . This variable  $T_i = 0$  if the ground truth is that the node is a normal region and  $T_i = 1$  if the ground truth is that the node is in an event region. We map the real output of the sensor into its binary bit-pattern, an abstract pattern variable  $P_i$ . This variable  $P_i = 0$  if the sensor measurement

indicates a change from the previous bit-pattern value and  $P_i = 1$  if it measures an unusual value.

There are thus four possible scenarios:  $P_i = \text{MSB change}$ ;  $M_i = 0$  (sensor correctly reports a normal reading),  $P_i = \text{IN-RANGE}$ ;  $M_i = \text{COSETS CODEBOOK}$  (when only 1 bit changes are reported by sensor reading),  $P_i = \text{IN-RANGE}$ ;  $M_i = \text{HUFFMAN CODEBOOK}$  (sensor correctly reports an unusual/event reading with majority reporting), and  $P_i = \text{LSB}$ ;  $M_i = 0$  (sensor reports very low readings). While each node is aware of the value of  $M_i$ , in the presence of a significant probability of a faulty reading, it can happen that  $M_i \neq T_i$ . We describe below a fault recognition algorithm to determine an estimate  $R_i$  of the true reading  $M_i$  after obtaining information about the sensor readings of neighboring sensors.

In our discussions, we will make one simplifying assumption: the sensor data is a geometrically distributed due to aggregating sensors which are placed in wireless range. Normal entropy is given for a random variable  $X$ , with  $n$ , outcomes  $x_i : i = 1, \dots, n$ , the Shannon entropy[12,13,14], a measure of uncertainty and denoted by  $H(X)$ , is defined as

$$H(X) = - \sum p(x_i) \log_b p(x_i) \quad (7.1)$$

where  $p(x_i)$ , is the probability mass function of outcome  $x_i$ , which is the average entropy.

One of a family of functionals for quantifying the diversity, uncertainty or randomness of a system. It is named after Alfrd Rnyi[12,13,14]. The Rnyi entropy of order  $\alpha$ , where  $\alpha \geq 0$ , is defined as

$$H_\alpha(X) = \frac{1}{1-\alpha} \log \left( \sum_{i=1}^n p_i^\alpha \right)$$

For any compression algorithm which assigns prefix codes and to uniquely be decodable. Let us define the **kraft Number** and is a measure of **the size** of  $L$ . We see that if  $L$  is 1,  $2^{-L}$  is .5. We know that we cannot have more than two  $L$ 's of .5.

If there are more that two  $L$ 's of .5, then  $K > 1$ . Similarly, we know  $L$  can be as large as we want. Thus,  $2^{-L}$  can be as small as we want, so  $K$  can be as small as we want.

Thus we can intuitively see that there must be a strict upper bound on  $K$ , and no lower bound. It turns out that a prefix-code only exists for the codes IF AND ONLY IF:

$$K \leq 1 \quad (7.2)$$

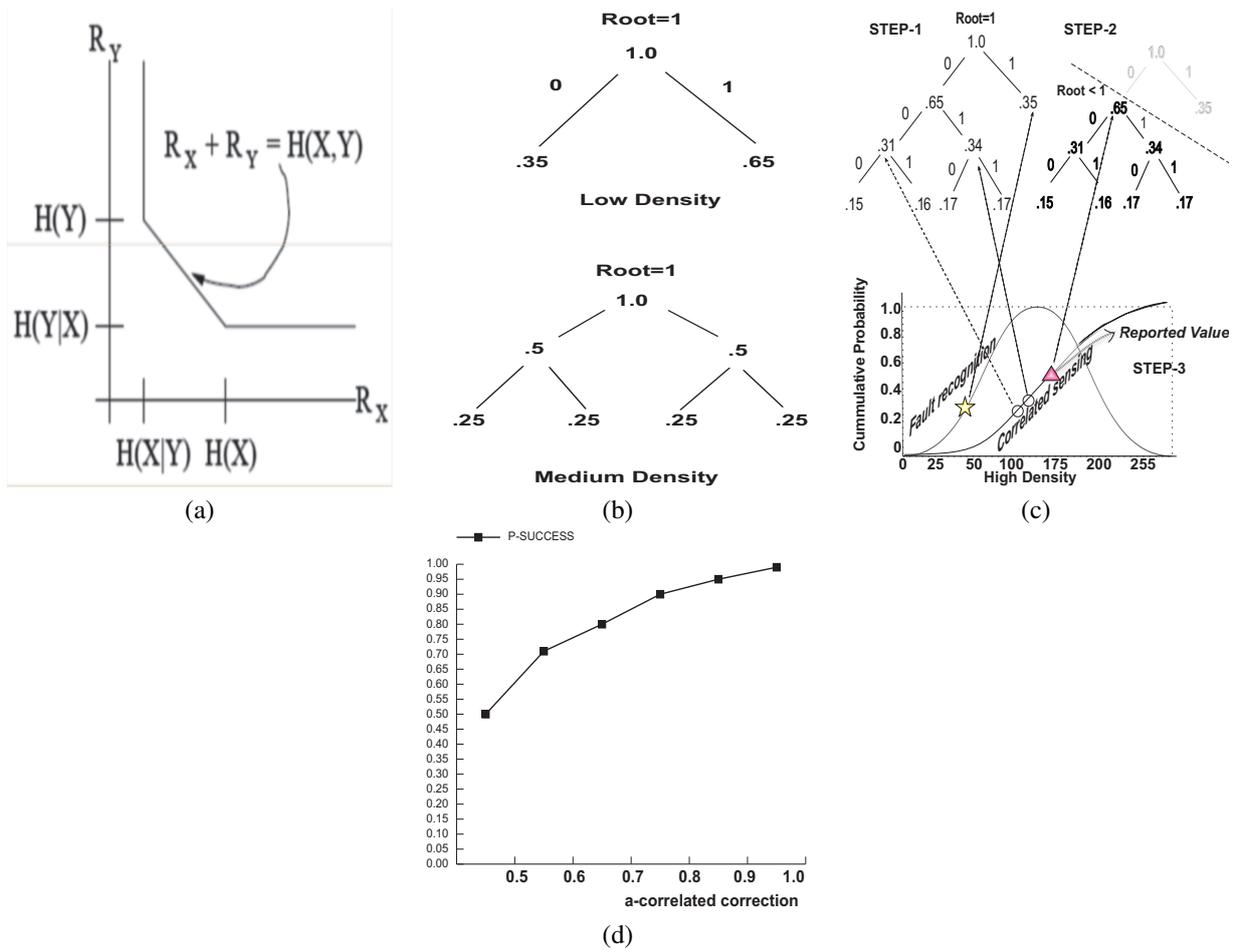


Fig. 6. (a) A graphical plot of the Slepian-Wolf rate (b) Huffman tree for small & medium (c) Faulty & corrected values (d) Success rate of Tx

The above equation is the Kraft inequality.

The sensor communication channel has a window of opportunity with a total duration (in bits) distributed geometrically with parameter  $a$ . The probability of successful transmission is given by probability mass function and lower-bound from equation 7.2  $0 < a < 1$ , find a code minimizing

$$\log \sum_{i \in \chi} p(i) a^{l(i)} \quad (7.3)$$

$$P_{success} = a^{L_a(p,l)} \quad (7.4)$$

From the above equation 7.4 which are used by compression algorithm to assign the prefix code optimally the quantitative sensor data optimization is due to  $a$  the geometrically distribution of values over time in  $P$  the probability [12,13,14] of success. The Huffman algorithm [12,13,14] combines the probability of the two least occurring values, here we will allow this step only if the probabilities differ and is  $< 2p(n-1)$ .

HUFFMAN ALGORITHM: *Modified priority Q for merging probability*

```

1: while !EndofPriorityQ do
2:   Select new item from Q
3:   if  $p(n) < 2p(n-1)$  then
4:     Merge  $Q = P(n) + P(n-1)$ 
5:   else
6:     Skip(n) Q
7:   end if
8: end while
    
```

VIII. ANALYSIS OF FAULT RECOGNITION ALGORITHMS

In order to simplify the analysis of the bitwise fault recognition mechanisms, we will make the assumption that, the network deployment can be categorized into low, medium and high densities. In the simple case compression algorithms cannot be able to use the measured statistics as there are few samples. A typical scenario is shown in figure 6(b). Typically the Huffman trees have few levels and can easily distinguish the correct measured value. As the density of the sensors increases the level of correlation and at the same time effecting the

range of measured values. As shown in figure 6(c) it has many levels and also many sensors typically have lower values or probabilities. For any compression algorithm there exist two components the encoder and the decoder. As discussed earlier in sensor networks due to high redundancy in measured data due to correlation. If the decoder reliably gets a reference value then the encoder needs to send only the difference or the change is bits, which reduces the number of bits needed to transmit. This is governed by Slepian-Wolf rate for the two correlated sources shown in figure 6(a). We further study the geometrically distributed property of sensors and how it factors into Huffman compression algorithm. We define a parameter  $a$  which is used in defining the codebook of the generated prefix code and also bounded by Kraft inequality defined in terms of individual code lengths. The success of transmission can be further calculated as shown plotted in Figure 6(d) by using the equation For a minimum pre-fix code  $a = 0.5$  as  $2^l \leq 1$  for a unique decidability

#### A. Iteration $a=0.5$

In order to extend this scenario to distributed source coding, we consider the case of separate encoders for each source,  $x_n$  and  $y_n$ . Each encoder operates without access to the other source. This scenario is illustrated in Figure 6. Using the techniques of fault recognition in this section, we would expect that each source could only be compressed to its marginal entropy. The anticipated efficiency of the algorithm is as follows in table II.

TABLE II  
ALGORITHMS

No	Algorithm	Compression Rate	Error Rate	Code Book
1	Slepian-Wolf	50%	Correctable Cosets	Decoder dependent
2	Huffman	70%	Single-bit (pre-fix)	Dic. based

#### B. Iteration $a \geq 0.5 \leq 1.0$

As in the previous case it uses correlated values as a dependency and constructs the codebook. The compression rate or efficiency is further enhanced by increasing the correlated cdf(as shown in the algorithm) higher than a  $> 0.5$ . This produces very efficient codebook and the design is independent of any decoder reference information. Due to this a success threshold is also predictable, if  $a = 0.5$  and the cost between  $L = 1.0$  and  $2.0$  the success = 50% and for  $a = 0.9$  and  $L = 1.1$ , the success = 71%. The efficiency of the algorithm is more adaptable and summarized in table II .

## IX. CONCLUSION

In this paper we look into pre-processing of sensor data streams to detect any outliers, which can generate false alarms. The local processing is done at the encoder, which has correlated information about the sensor data, which helps in reducing redundant transmission. The encoding algorithm uses a coefficient, which is independent of the measured value

but consistent with amount of redundancy in the data. This coefficient is learnt overtime and the sensors are allowed to calibrate by the assumption that the variations in data would be around the learnt correlated coefficient.. The simulation shows that efficiency for a sensor network is high due to correlated sources as shown in algorithm1. Algorithm2 uses standard Huffman code by adding a success rate probability, which helps to reduce energy and the code length over successive transmissions.

## ACKNOWLEDGMENT

The authors would like to thank Louisiana State University Department of Computer Science which helped in supporting the preparation of this writing which in turn helped the defence of the PhD thesis and candidature. This work was supported by the Post Katrina funds (PKSFI) from Board of Regents and also from DoD-DEPSCoR-Office of Naval Research Program at LSU, Baton Rouge.

## REFERENCES

- [1] Iyer, V.; Iyengar, S. and Rammurthy, G. Distributed Source Coding for Sensor Data Model and Estimation of Cluster Head Errors using Bayesian and K-Near Neighborhood Classifiers in Deployment of Dense Wireless Sensor Networks SENSORCOMM, 2009.
- [2] Iyer, V., Garimella, R. M., and Srinivas, M. B. 2008. Min Loading Max Reusability Fusion Classifiers for Sensor Data Model. Second international Conference on Sensor Technologies and Applications - Volume 00 (August 25 - 31, SENSORCOMM 2008).
- [3] Bhaskar Krishnamachari, Member, IEEE, and Sitharama Iyengar, Fellow, IEEE., Distributed Bayesian Algorithms for Fault-Tolerant Event Region Detection in Wireless Sensor Networks, IEEE Transaction On Computers, Vol. 53, NO. 3, March 2004.
- [4] Vasanth Iyer, S.S. Iyengar, N. Balakrishnan, Vir. Phoha, M.B. Srinivas. FARMS: Fusionable Ambient Renewable MACS. Sensor Application Symposium (SAS-2009), IEEE 9781-4244-2787, 17th-19th Feb, New Orleans, USA.
- [5] Training Data Compression Algorithms and Reliability in Large Wireless Sensor Networks. Vasanth Iyer, Garimella Rammurthy and M.B. Srinivas. International Journal On Smart Sensing and Intelligent Systems, Vol. 1, No. 4, Dec.'08 ISSN 1178-5608.
- [6] Environmental measurement OS for a tiny CRF-STACK used in Wireless Network Vasanth Iyer, G.Rama Murthy and M.B. Srinivas- Sensors & Transducers Journal-April 2008, ISSN 1726-5479 2006 by IFSA.
- [7] Richard O. Duda, Peter E. Hart, David G. Stork, Pattern Classification. John Wiley, 2005.
- [8] Wendi Heinzelman, Hari Balakrishnan, Anantha Chandrakasan, LEACH (Low Energy Adaptive Clustering Hierarchy) MIT,
- [9] Vasanth Iyer, S.S. Iyengar, G. Rama Murthy and M.B. Srinivas. 'Multi-hop Sleep Scheduling and Local Data Link Aggregation Dependant Qos in Modeling and simulation of Power-aware Wireless Sensor Networks IWCMC 2009 Wireless Sensor Networks Symposium
- [10] Vasanth Iyer, G.Rama Murthy and M.B. Srinivas. Entropy based Variable Rate compression for Low-bandwidth Multi-Media Streams. DFMa, Oct, 2008, Penang, Malaysia.