

Statistical Model of Lossy Links in Wireless Sensor Networks

Alberto Cerpa, Jennifer L. Wong, Louane Kuang, Miodrag Potkonjak and Deborah Estrin
Computer Science Department, University of California, Los Angeles, CA 90095
{*cerpa,jwong,lkuang,miodrag,destrin*}@cs.ucla.edu

Abstract—Recently, several wireless sensor network studies demonstrated large discrepancies between experimentally observed communication properties and properties produced by widely used simulation models. Our first goal is to provide sound foundations for conclusions drawn from these studies by extracting relationships between location (e.g. distance) and communication properties (e.g. reception rate) using non-parametric statistical techniques. The objective is to provide a probability density function that completely characterizes the relationship. Furthermore, we study individual link properties and their correlation with respect to common transmitters, receivers and geometrical location.

The second objective is to develop a series of wireless network models that produce networks of arbitrary sizes with realistic properties. We use an iterative improvement-based optimization procedure to generate network instances that are statistically similar to empirically observed networks. We evaluate the accuracy of our conclusions using our models on a set of standard communication tasks, like connectivity maintenance and routing.

I. INTRODUCTION

It is well known that the performance of many protocols and localized algorithms for wireless multi-hop sensor networks greatly depend on the underlying communication channel. Hence, to evaluate performance in simulation, we must have an accurate communication model. Until recently, only two approaches have been in widespread use in the sensor network community: unit disk modelling and empirical data traces.

Both approaches present some drawbacks. For example, the unit disk model implies complete correlation between the properties of geometric space and the topology of the network, a property refuted by numerous experiments in actual deployments [1], [2], [3]. When using empirical data traces approach is difficult and expensive to create a large number of large networks that are properly characterized. Therefore, neither probabilistic nor statistical analysis of large networks is feasible. In addition, since a given trace is the result of communication over a specific topology, such a trace does not permit a simulator to reposition nodes. Finally, without validated communication analysis, theoretical analysis is not possible.

In an effort to address this problem, recently there have been a number of efforts to empirically capture communication patterns in wireless sensor networks. In particular, there have been several studies that use different low power, narrow band radio transceivers chipsets [4], [5] to deduce properties of communication links in wireless networks in several environments, such as open space and laboratories. These hybrid models introduce empirically observed factors that modify the communication patterns based on the unit disk communication model.

While these models are a significant step forward with respect to the unit disk model, they are only an initial step in the exploration of the space. These initial models do not capture many important features of communication links in empirically observed networks. For example, they do not address the correlation in communication reception rates between nodes that originated at the same transmitter or differences in the quality of transmitters.

Our goal is to develop accurate simulations of sensor network communication environments that are statistically accurate with respect to

several features that impact network protocols and algorithms in real networks. To generate these simulated environments, we construct a set of models that map communication properties such as absolute physical location, relative physical proximity and radio transmission power into probability density functions describing packet reception likelihood. For all of these models, we calculate an interval of confidence. These models not only serve to generate simulated environments, they themselves have lent support to many hypotheses relating to variation in communication link quality [1], [2]. In our study, we do not consider packet losses introduced by multi-user interference (concurrent traffic, contention-based MAC). Nevertheless, our results are useful for three reasons. First, the amount of traffic expected in most application in sensor networks is small, which means either small contention, or in case of highly synchronized events, nodes could be programmed to prevent simultaneous transmissions. Second, they apply directly when using contention free MAC protocols, like pure TDMA or pseudo-TDMA schemes [6]. Finally, they provide a tight upper bound as to what is achievable when using contention-based MAC schemes. The analysis of multi-user interference and temporal properties of the links is part of future work.

II. RELATED WORK

There is a large body of literature on mobile radio propagation models that have influenced this work. The emphasis has been on large scale path loss models that predict the average received signal strength at a given distance from the transmitter and the variability of the signal strength in proximity to a particular location [7]. Furthermore, the models are used to predict the coverage area of a transmitter. In addition, small scale fading models are used for modeling localized time durations (a few microseconds) and space locations (usually one meter) changes. All these models are based on the Fries free space equation and indicate that reception quality decays with the inverse of distance raised to a small power [7].

Differences between the classical models and our approach are numerous. We have different modeling objectives (reception rate of packets vs. signal strength), our radios have different features (e.g. communication range in meters instead of km), we capture phenomena that is not addressed by the classical channel models (asymmetry, different quality of receivers and transmitters, correlations between reception rate of links), we use different modeling techniques (free of assumptions), and we use unique evaluation techniques (resubstitution and evaluation of multi-hop routing). We believe existing and new techniques have complementary objectives, tools, and applications.

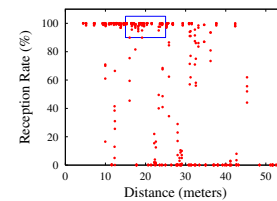
More recently there have been many studies of significant-scale deployments in several environments [2], [8], [3], [1], [9], [10]. Majority of these studies used the TR1000 [5] and CC1100 [4] low power RF transceivers. There are three major differences between the models developed in this paper and the previously published models. The first is that we study the impact of a significantly large number of factors that impact reception rate and attempt to model not only isolated pairs of transmitters and receivers, but also the correlation between different pairs and different subsets of links. The

```

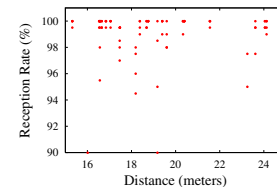
1. Conduct exploratory data analysis;
2. while (interval of confidence > criteria) {
3.   Collect new data or define new windows;
4.   Sort all points according to distance;
5.   for(from smallest to largest distance) {
6.     Define sliding window for distance;
7.     Apply weight function to distances inside of sliding window;
8.     Sort all points according reception rate;
9.     for(from smallest to largest reception rate ) {
10.      Define sliding window for reception rate;
11.      Apply weight function to reception rates sliding window; } }
12. Build mapping function;
13. Build normalized mapping function;
14. Establish intervals of confidence; }

```

Fig. 1. Pseudo-code of the PDF model generation for two features.



(a) Data set



(b) Zoom data box (a)

Fig. 2. Scatter plot of distance vs. reception rate

second major difference is that our goal is not only to establish a model, but also to establish statistically sound measures as to what extent the model corresponds to experimentally captured data. Our statistical techniques are generic in the sense that they can be used in other studies with minimal changes. Also, we have developed a procedure for creating instances of an arbitrary size for simulation and mechanisms to ensure that they are accurate models with respect to the collected data.

All of our techniques are based on non-parametric statistical procedures. Specifically, we use smoothing and density kernel estimators, resubstitution and bootstrap techniques [11].

III. INDIVIDUAL LINK MODELS

In this Section, we present a new statistical approach for building communication link models in wireless multi-hop networks. The goal is to find a statistically sound mapping between two user-specified features that characterize communication links.

A. Design Guidelines

Our starting task is to analyze the dependency between two properties of wireless networks. We note that exactly the same procedure described below can be used to find the dependency between *any* two wireless communication features, but for the sake of brevity and clarity, we focus on two specific features: distance and reception rate. The objective is to find the PDF of reception rate for any distance and to calculate intervals of confidence. For example, we could use our model to find that the probability of the reception rate of the link to be 95% at a distance of 25 ft is 0.05 ± 0.000012 .

We are guided by three principles, smoothness, compactness, and prediction ability. The validation for adopting these principles is provided by evaluation using resubstitution, which indicates that the derived models have tight interval of confidence and therefore, the statistical model is accurate and the assumptions are justified. The smoothness property states that if two pairs of receivers have very similar distance, their reception rates also often have rather similar probabilistic distribution. In other words, instances of reception rates may be different from one distance to the other, but the underlying reception rate probabilistic distribution is similar. There are two fundamental justifications for this assumption. The first is that at

an intuitive level one expects that small changes in one variable (in our case, distance) should have limited impact on the probability distribution of the other parameter (reception rate). In addition, it is important to recognize that both distance and reception rate are subject to errors in measurement that smooth the mapping function.

Quality of the statistical model is ensured through compactness and performance on test cases to measure the prediction ability. There are two sound criteria for any sound statistical model. The first is the Occam principle: the ability to explain a large set of data using a small number of parameters is usually a strong indication that the model will predict well. From a statistical point of view, our goal is to simultaneously have low bias and low variance and therefore low prediction error. Low bias is ensured by preferring models that use fewer parameters. For this task we use Akaike information criteria (AIC) [11]. The second criteria is its ability to predict. We scan and alter various parameters in our procedures so long as the adopted parameter values produce a model that withstands standard evaluations of accuracy. Specifically, we use the resubstitution rule, which builds additional models using a variety of randomly selected subsets of the data set. If the resulting models from all data subsets are similar, we conclude that the parameters used were properly selected.

From a technical point of view, when building a model of individual links we have two major difficulties: (i) we do not have enough measurements for each distance of interest, and (ii) for a given distance and given reception rate, we do not have enough collected data samples. We use the kernel smoothing technique to resolve this, and identified that the best performing window had $\pm 10\%$ of the size of the central value and pyramidal shape.

B. Methodology

The global flow of our approach is shown in Fig. 1 using pseudo-code format. The starting point for the procedure is exploratory data analysis. As the first step of this phase, we examine a scatter plot of all available data points. Specifically, we position each communication link in a two dimensional space, placing on the x-axis distance and on the y-axis the reception rate. The goal is to identify if there are any specific trends in the data and to determine whether it is advantageous to split the data into two or more subsets that have specific features. Fig. 2(a) illustrates a scatter plot of distance versus reception rate

TABLE I
GLOBAL EVALUATION RESULTS

Environment	Conf. Level	H.L.PDF Value	Conf. Intervals
Indoor	90	0.997627	± 0.325969
Outdoor	90	1.064365	± 0.381719
Indoor	95	1.023886	± 0.723887
Outdoor	95	1.022372	± 0.691752

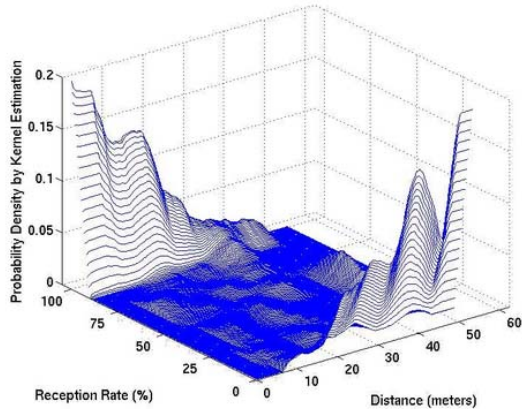


Fig. 3. PDF for distance versus reception rate.

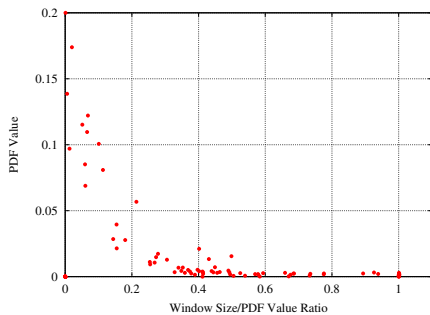


Fig. 4. PDF values of the different random points as a function of the confidence interval/PDF value ratio. Outdoor Urban, 90% Confidence Level.

at medium power for the outdoor case. Fig. 2(b) shows a zoomed version of a subset of data which was examined during exploratory data analysis. We conducted exploratory analysis in order to gain semantic insights, which can only be done by humans according to the techniques presented in [12]. We did consider automatic clustering, in particular, self-organizing maps, principal components, independent components and multidimensional scaling [11], but they did not show intuitive trends.

Phase two consists of three steps shown in Fig. 1 in lines 4-8. In the step shown in line 4, we sort all available data according to distance to identify data points that are similar with respect to this parameter. Next, we use a sliding window for all points which are within a similarity range of a given point (distance). Each point within this range is weighted according to its quantified similarity to a given point. For each distance of interest we also build another system of sliding windows this time along the y-axis corresponding to the reception rate. The points within the window are weighted as the product of the weight factor of both the distance window and the reception rate window.

Once the first eleven lines of the pseudo-code are executed we have enough information to build a PDF that indicates how likely a particular reception rate is for a given distance. For this task, following compactness principle, we used quadratic least linear squares fitting for a particular pair of distance and reception rate.

Once the model is built, the next step is PDF normalization that ensures that for a given distance the integral of the function below

the PDF mapping function is equal to one. Fig. 3 illustrates how the normalized reception rate PDF changes with respect to distance.

C. Evaluation

The final step of our procedure is the evaluation of the quality of the developed statistical model for the PDF. The evaluation procedure itself has three components: Monte Carlo sampling, resubstitution, and establishment of interval of confidence. Monte Carlo sampling selects k (in our experimentation we use 200) randomly selected pairs of distance and reception rate points.

Resubstitution is the process where a statistical model is built using the exact same procedure (same kernel window scope and weight function) on randomly selected subsets of data. Specifically, in our simulations, we select 70% of the available data to build a model on each resubstitution run. For each resubstitution run we record the value of the PDF function at each of the k selected points. After conducting m resubstitution runs (in our experimentation m was 100), we are ready to establish an interval of confidence for our statistical PDF model. This is performed in two stages. We first establish an interval of confidence for each point individually, and then by combining information from all local interval of confidence we establish a global interval of confidence. Fig. 4 shows the relationship between the different confidence intervals for each random sample tuple (reception rate and distance) and the highest likelihood PDF value for different confidence levels. Each point in the graphs show the highest likelihood PDF value with its confidence interval. For example, the top left point in the graph of Fig. 4 corresponds to sample point of distance 52 meters, reception rate 0% with highest likelihood PDF value of 0.2 ± 0.0001953 with confidence level of 90%. The final step of resubstitution is to build a global measure of the model's accuracy. To build a global interval of confidence we use the following procedure. First, for each separate point in k , we use the highest likelihood PDF value and normalize all other values against this value. After that, we combine all data from all sampling points into one set of the size $k \times m$. Finally, we calculate the confidence intervals of the normalized global array. Table I shows the overall interval of confidence for indoors and outdoors with different confidence levels. In general, the global highest likelihood PDF values are centered around one, which is a good sign of the statistical soundness of the model.

IV. EXPERIMENTAL DATA COLLECTION

We used an existing data set and performed additional experiments using the SCALE wireless measuring tool [1]. The topology used for our experiments consisted of 16 nodes distributed in an ad-hoc manner in different environments. We also used up to 55 nodes for our indoor experiments deployed in the ceiling of our lab. For outdoor experiments, nodes were placed in a variety of different positions, such as near the ground or elevated off the ground, with or without line of sight (LOS) between them, and with different levels of obstructions (furniture, walls, trees, etc.). The placement of the nodes also took into account the distance between them, in order to

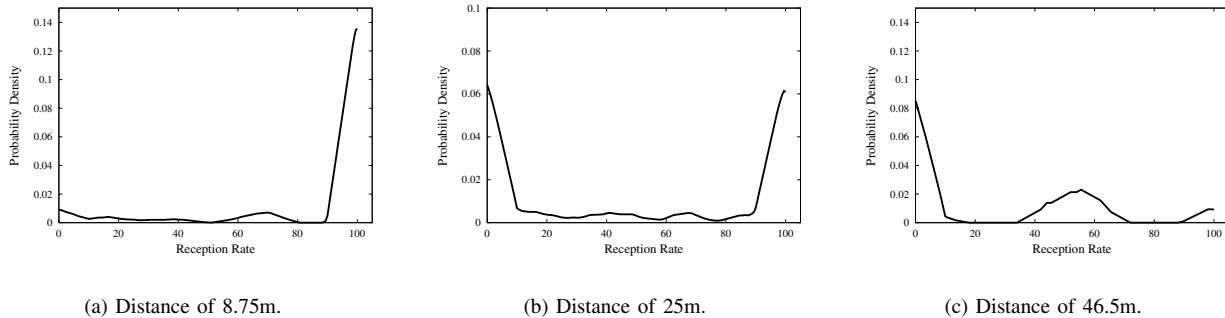


Fig. 5. PDFs for reception rate for various values of distance.

TABLE II
SUMMARY OF EXPERIMENTAL FACTORS CHANGED

Environment	Will Rogers State Park (outdoor), Boelter Hall courtyard (urban), LECS lab (indoor)
Radio Type	TR1000 [5] (916MHz), CC1100 [4] (433MHz)
TX Power	-10 to 0 dBm (TR1000), -20 to 10 dBm (CC1100)
Packet Size	25, 50, 100, 150, 200 bytes
Antenna Height	0, 0.25, 1, 3, 5 ft

create a rich set of links at distances varying from 2 to 50 meters and in multiple different directions from any particular sender.¹ In most of our experiments, each node sends up to 200 packets per round, transmitting 2 packets per second.² Using this setup, we varied several factors in our experiments. Table II shows a summary of them.

In summary, the data set used in this paper consisted of packet delivery data from more than 450,000 packet probes in experiments performed in 3 different environments, with 2 different type of radios, with 6 different power settings, 5 different packet sizes, and 4 different antennae heights.³ We used up to 16 nodes in our outdoor experiments and up to 55 nodes in our indoor experiments. We measured the packet delivery performance of 240 links for the outdoor experiments and 2970 links for the indoors experiments.

V. PROPERTIES OF INDIVIDUAL LINKS

At the highest level of consideration the features can be classified into two groups: physical properties of the network, and communication features of the network. Physical properties include distance, direction as a function of angle with respect to reference direction, and areas. Communication properties include reception rate between receiver A and transmitter B, asymmetry in communication which refers to the absolute difference in reception rates between a pair of nodes (transmitter A \rightarrow receiver B and transmitter B \rightarrow receiver A), and temporal variation of reception rate between receiver A and transmitter B.

We have analyzed two types of mapping functions between properties. The first is the established one-way mapping relationship between a given structural property and a targeted communication

¹There are algorithms to find the optimal placement of nodes to get a uniform range of distances in the area of interest. We did not perform that optimization in our experiments.

²We have left for future work the evaluation of how accurate is an average reception rate.

³We did not test all possible combinations in all environments.

feature. The second analyzes the one-way dependency between two communication features. We once again emphasize that our goal is to not only establish the most likely value of one property for a given value of another property, but also to obtain probability distribution functions for all expected values of the second property for a given value of the first feature. We have studied the following pairs of properties.

Dependency of reception rate as a function of distance. This property is selected mainly because there is a wide consensus that distance significantly impacts reception rate.

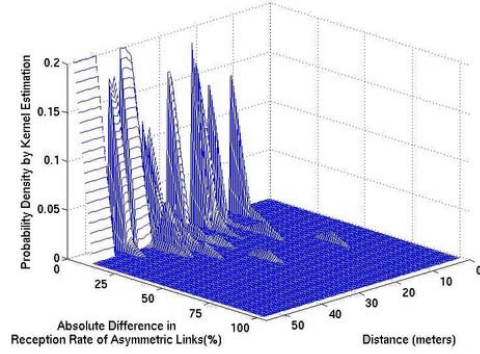
Dependency of asymmetric reception rate as a function of distance. Note that in the previous case we assumed that reception rate between transmitter A to receiver B is the same as from transmitter B to receiver A, but recently several empirical studies demonstrated this is not the case [1], [2]. Our goal is to quantitatively capture how frequently there is asymmetry in reception rates as a function of distance.

Dependency of asymmetric reception rate as a function of reception rate. This property studies functional dependencies between two communication properties. Our goal is to identify if it is more likely that high asymmetry happens when links have high, low, or medium reception rates. For example, we are interested if it is more likely to have a pair of nodes with reception rates of 95% and 75%, or with 30% and 10% reception rates.

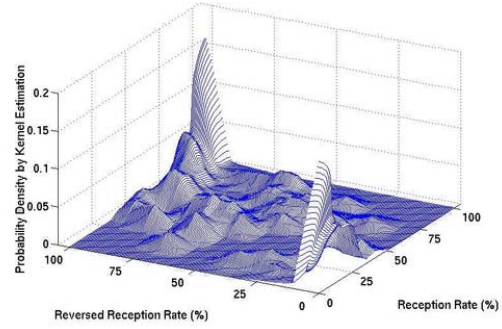
Dependency of reception rate standard deviation as a function of the average reception rate. The final property studies temporal dependencies between two communication properties. An empirical study [1] has shown that such correlation exists. Our goal is to quantitatively capture this relationship and provide some initial results on how this property affects the link estimation algorithms used for online quality estimation.

In addition to the listed properties, we also studied link quality dependency on angle, but were not able to identify any interesting patterns with significantly strong intervals of confidence.

In Sect. III, Fig. 3 we have illustrated how the reception rate changes as a function of distance. Figs. 5 show normalized PDFs for three typical distances for 8.75, 25, and 46.5 meters. These results confirm the findings of several studies in the literature that show that there is a significant percentage of the radio range where links are highly variable, with similar probabilities of having very high or low reception rates. In addition, we show that even for very short distances, the probability of having very low reception rate links is not zero, and it starts growing fast as distance increases. More importantly, it is clear from the graph in Fig. 5(b) that the average and standard deviation values of reception rate are insufficient parameters



(a) Asymmetric links vs distance



(b) Asymmetric links vs reception rate

Fig. 6. PDFs for asymmetric links features.

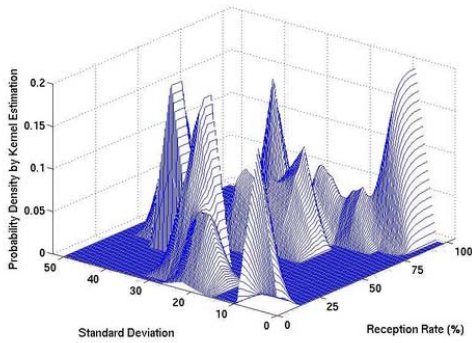
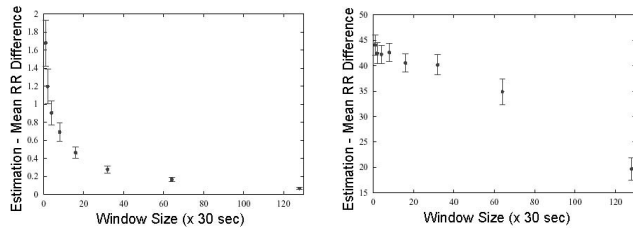


Fig. 7. PDF for temporal variation as a function of the reception rate.



(a) High Reception Rate

(b) Med. Reception Rate

Fig. 8. Time series for on-line link quality estimation.

to model reception rate as a function of distance. While the average reception rate is around 50% in this case, most of the links have either very high or low reception rates.

Figs. 6(a) and 6(b) show the PDF of how asymmetric reception rate depends on distance and average reception rate. Fig. 6(a) shows that there is no clear correlation between link asymmetries and distance. Fig. 6(b) shows an interesting pattern; links with very high or very low reception rates tend to be highly symmetrical, as it can be observed by the two peaks in the PDFs. Links with medium reception rate tend to be much less symmetrical.

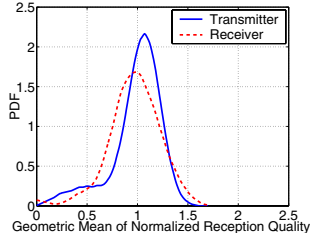
Fig. 7 shows the temporal variability of the links as a function of

the reception rate. We clearly see that links with very low or very high reception rates tend to be more stable over time (smaller standard deviation), while the links with intermediate values of reception rate tend to be more unstable (higher standard deviation).

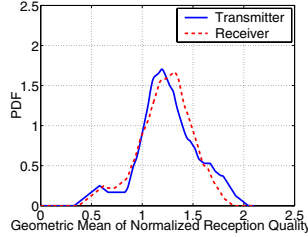
From our data, we observe that while the quantitative values of the PDFs for different conditions were not the same, the PDFs generated were qualitatively similar in most cases. For example, the PDF shown in Fig. 3 was qualitatively similar across all three environments tested. We have left the statistical analysis of the differences between the different conditions to get statistical sound conclusions for future work because it requires additional experiments.

One interesting question we wanted to answer is how long a node needs to measure the communication channel in order to get an accurate estimate of reception rate with a certain confidence interval. This has a profound impact in the design of algorithms for topology control that need to measure the channel as little as possible in order to save energy by periodically turning the radio off. To evaluate this, we took long time series of reception rate data, and picked k window sizes. For each window size, we took p (set to 100) initial random points of measurements from the time series, generating a reception rate estimate for each p using only a window of size k (ranging from 30 seconds to 64 minutes) of data from the starting point. Then we compare the absolute difference between each of the $p \times k$ estimates with the absolute reception rate calculated using the entire time series of data. Fig. 8 shows the results of the previous analysis on two qualitatively different type of links. Fig. 8(a) shows that links with very high reception rate need very short window sizes to get an accurate estimate of the reception rate, and they converge quite fast to an accurate estimate (low reception rate links show similar behavior). Fig. 8(b) shows that links with intermediate reception rates take much larger window sizes to converge to accurate estimate values. We have left for future work the issue of optimal on-line link characterization using statistical methods.

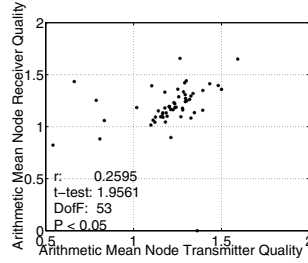
From the spatial, asymmetrical, and temporal properties presented in Figs. 3, 6(b) and 7 we can see an interesting pattern that has emerged. For a large range of distances there is a low but non-zero probability of links with medium reception rates. These links are also the ones that present the most highly asymmetrical and temporal variability properties. We believe these links may introduce serious stability and convergence problems for several routing algorithms, and it might be useful to design mechanisms to detect these types



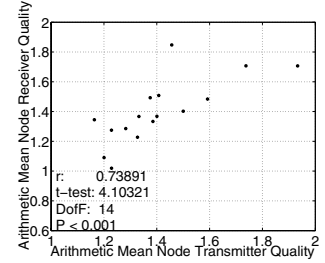
(a) Mean Quality of Reception: Indoor.



(b) Mean Quality of Reception: Outdoor.



(c) Mean of Quality Correlation (per node): Outdoor.



(d) Mean of Quality Correlation (per node): Indoor.

Fig. 9. PDF for Normalized Transmitter and Receiver Quality and Correlation

of links and filter them out [9]. Another interesting point is that reliable, highly symmetrical and stable links exist even at long distances (although with low probability). It is important to detect and take advantage of those long distance/high quality links in order to minimize packet transmission in a multihop setting. Online detection and use of these type of links could affect algorithm design. (e.g. by minimizing energy consumption or end-to-end hop count).

VI. GROUP LINK PROPERTIES

Group link properties are joint properties of related subsets of links. These links may be links that originate from the same transmitter or received by the same receiver, processed by the same radio, or communicated by nodes that are geometrically close. These properties are of crucial importance for any analysis and answer the frequently asked fundamental questions about reasons for particular behavior of communication patterns. These questions include whether the performance of a particular node as a transmitter mainly depends on the quality of its radio or its geometric position. Another frequently asked question is whether asymmetry is a consequence of different radio properties between two nodes or their location. However, with the exception of the property which examines pairs of links between two nodes, group link properties have been rarely studied due to their perceived complexity.

The first question we answer is to what extent the quality of transmitters and receivers on different nodes is uniform. We normalized the quality of each link at each node versus the average link quality at the corresponding distance in terms of reception rate. After that, we calculated the geometric mean of all links that originate or end at a particular node. For the reception quality, we decided to use the geometric mean instead of arithmetic in order to avoid too high an impact from a few exceptionally strong links. For example, if a node has one link that is 5 times better than average and 5 links that are 5 times worse than average, the arithmetic mean would still indicate that the nodes have links of superior quality, which is obviously not the case. When there are not outliers, the arithmetic mean is preferred.

Figs. 9(a) and 9(b) show the PDFs for normalized transmitter and receiver quality of nodes in indoor and outdoor environments. We see that a large percentage of nodes are either significantly better transmitters or receivers than average, in particular for the outdoor environment. The result is important because if nodes were completely uniform, deployment strategy would be based solely on topology. However, our results show that further communication optimization can be made by considering transmitting/receiving quality.

The second question we answer is whether there is correlation between the quality of the transmitter and receiver on the same node.

TABLE III
CORRELATION OF ALL PAIRS FOR INDOOR AND OUTDOOR.

	Outdoor			Indoor		
	r	t-test	DoF	r	t-test	DoF
TX	-0.004	0.121	887	0.0592	11.824	39770
RX	0.012	0.370	885	0.0590	11.859	40183

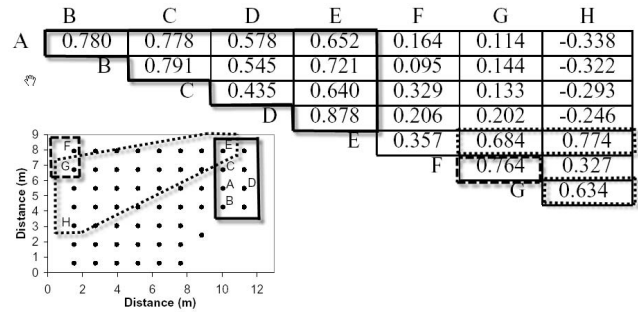


Fig. 10. Covariance Matrix and Layout for Indoor Experiments.

Figs. 9(c) and 9(d) show a summary of our results. We analyzed both indoor and outdoor data using arithmetic mean. We calculated both arithmetic and geometric mean correlations, but due to the lack of outliers in the data, we preferred to use the arithmetic mean. All studies indicate that there is a positive correlation of transmitting and receiving capability of the nodes, and the probability of this result being accidental is low (lower than 0.1% in the indoor case). The linear correlation factor values are different depending on the environment, being much higher for the indoor case.

Once we conclude that some nodes are much better transmitters or receivers than other nodes, the natural question is to what extent they are uniformly better transmitters or receivers with respect to all their links. In order to answer this question we calculated the correlation between all transmitting (receiving) links related to the same node. Table III shows the correlation value r , the t-test value and the degree of freedom (DoF). For both indoor and outdoor environments we see essentially very small or no correlation with very high probability (the probability of this result being accidental is lower than 0.1% for the indoor case). This essentially means that no node has perfectly good links to all other nodes in some distance range, and even the best nodes have average or very poor links. In addition, almost all nodes have good links to some neighbor in the same distance range.

The last question we would like to answer is whether there

are subsets of nodes that communicate well with each other while communicating at significantly lower levels with other nodes in the network. Fig. 10 shows the covariance matrix for 9 nodes in the indoor environment. We clearly see that nodes A, B, C, D and E form one group, nodes E, G and H another, and F and G, the third group. All nodes in these groups are highly correlated in terms of normalized communication with respect to other nodes. The data was obtained in the following way. For each node we sorted in decreasing order the quality of its links to other nodes. After that, for each pair of nodes, we found a subset of corresponding receivers that hear both nodes, and eventually found rank correlation for these two lists. As part of the table indicates, very often the correlation between two nodes is rather high, close to positive 1 or very low close to -1. The Spearman t-test indicates that all covariance values have probability of accidentally happen well below 0.1%. In other words, group of nodes in a particular distance range can communicate to each other *significantly* better than other group of nodes in the *same* distance range. Identification of these groups of nodes could be important for tree-based routing algorithms; it would be convenient that at least one node in each of these groups join the tree since it could communicate better to the other nodes in the group than any other node.

VII. WIRELESS NETWORKS GENERATORS

Using the knowledge gained from analysis of single and multiple link properties, we have built a series of wireless multi-hop network instance generators to be used in simulation environments. We present three models, increasing in complexity, which create communication links for an arbitrary network that are statistically similar to observed networks. The basic model assigns communication links based solely on the relationship between reception rate and distance. To build the more complex models, we introduce an iterative improvement-based procedure for creating communication links which abide by multiple link properties. The starting point for all models is the generation of a user specified number of nodes in the given area, with specific locations or a particular distribution.

A. Probabilistic Disk

The basic model, probabilistic disk, considers only the dependency between distance and reception rate. It is created by generating for each calculated distance between two nodes a randomly selected reception rate according to the PDF (Fig. 3) for the respected distance. We first translate the PDF into the corresponding cumulative distribution function (CDF) and use a uniform random generator between 0 and 1 to generate a value of CDF. The resulting value is then mapped into the corresponding reception rate.

B. Bi-Directional Correlated Probabilistic Disk Model

In this model, we consider two functional dependencies. In addition to the dependency between reception rate and distance (Fig. 3), we also consider the dependency between asymmetric reception rate and reception rate (Fig. 6(b)). Our goal is to generate an instance of the network where all communication links follow the PDF for a corresponding distance *and* for any given pair of nodes, we have a reception rate between transmitter A and receiver B and transmitter B and receiver A that follows the PDF for asymmetric reception rate. An instance of the network which follows this model is generated in the following way. We first generate for each pair of nodes the reception rate between the transmitter of node A and receiver of node B, where the notation of nodes A and B for a given pair of nodes is randomly conducted. Next, we generate the reception rate of the transmitter of node B and the receiver of node A following

the PDF for reception rate into probabilistically selected asymmetric rate using one of the previously mentioned methods. One can prove using Bayesian rules that the network generated using this procedure does follow both PDF functions.

C. Non-parametric Statistical Model

While the Bayesian rule is powerful enough to generate instances of the network that follow one and in some cases two PDFs, it is easy to see that when a large number of statistical measures must be followed, it does not provide an adequate solution. For example, it is not clear how to simultaneously generate a network which follows PDFs for reception rate versus distance, asymmetric reception versus distance, and non-uniform quality of transmitters and receivers. In order to overcome this difficulty we have developed an iterative improvement-based algorithm that generates an instance of the network that approximately, or arbitrarily closely, follows an arbitrary number of interacting PDFs defined on arbitrary pairs of network and communication properties. The key idea is to first separately generate an instance of the network that follows each of the considered PDFs and to randomly select one of them as a starting point for the iterative improvement procedure. At each step of the iterative improvement procedure, we attempt all possible changes at all possible pairs of nodes A and B and select one which makes the overall discrepancy between the parameters of that network more similar to a combination of the originally generated networks that separately considers the PDF of only a single property. The similar PDF function is defined using standard L_2 measure. The procedure is repeated until no further improvement can be found. In order to improve the quality of the generated network, one can perform restarts or employ probabilistic mechanisms for escaping local minima (e.g. simulated annealing). If a restart is performed, there is the option to probabilistically select one of the final solutions for the restart according to their maximum likelihood expectation. These expectations are generated from the space that contains all networks that follow all the specified PDFs.

D. Generation of Large Network Instances

Scalability is one of the key issues in wireless sensor networks both during deployment as well as during protocol and algorithm development. Unfortunately, it is both expensive and time consuming to deploy large networks solely for the purpose of building a model or developing a localized protocol. Therefore, there is a need to develop a methodology and approach that creates and validates networks of an arbitrary size.

We have developed a perturbation-based analysis that facilitates sound statistical validation of networks of an arbitrary size with respect to experimentally available and characterized networks. The key idea is to begin by creating an instance of the network using a specific communication model, and run the algorithm or protocol of interest on the network instance. Next, we replace randomly selected subparts of the instance with instances of data from actually deployed networks. generate, using the selected statistical model, the connectivity between nodes in the patches of the real networks and the neighboring nodes from the generator. After the procedure is completed, we compare the initial and perturbed networks with respect to results they produce on a task of interests. The extent to which the results are similar, the large instance is representative of the real-life networks. Fig. 11 illustrates the similarity in terms of all-pairs shortest path between two large network examples of 400 nodes built using the asymmetric link model.

TABLE IV
COMPARISON OF FOUR STATISTICAL MODELS USING FLOYD-WARSHALL ALL PAIR SHORTEST PATH ALGORITHM.

	Unit	Unit Real	Prob.	Prob-Real	Asymm	Asymm Real	Statistical	Statistical-Real
MIN	2	2	2.0079	2.0079	2.00188	2.00188	2.00002	2.00002
MAX	26	20.0569	41.881	43.354	45.9964	44.1535	42.99	42.9285
AVE	6.87574	5.78918	14.687	15.002	14.8176	14.6217	14.6991	14.6928

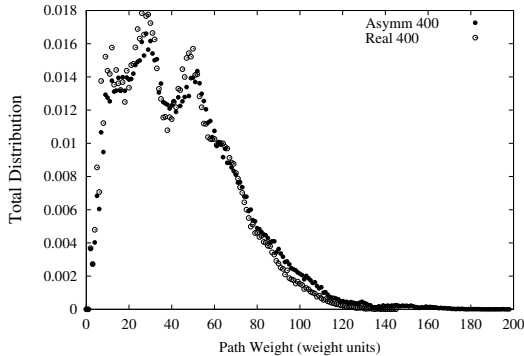


Fig. 11. Similarity between path weights in large networks.

We compared using the perturbation-based method four models: unit disk model, probabilistic disk model, asymmetric probabilistic disk model, and a non-parametric statistical model. For this purpose we compare the length of all-pairs shortest paths for an instance with 400 nodes. Table IV provides a summary for the length of the minimal, maximal, and average path. Note, that all three newly developed models, and in particular the non-parametric statistical one, are much more statistically sound.

The network generators have been implemented in the Em-Star simulator [13]. The code is available for downloading at: <http://cvs.cens.ucla.edu/emstar>.

VIII. DESIGN CONSIDERATIONS AND CONCLUSION

From the conceptual point of view, the first important observation is that the distribution of lossy links can greatly affect routing algorithms based on geometric concepts. For example, all local avoidance approaches that reduce the routing problem to traversal on Gabriel or local neighborhood graphs may no longer be applicable in practice. Another, possibly more impacting ramification is that no deterministic method can be used to guarantee packet delivery in stateless routing protocols. This is justified by the small but non-zero probability of having links with *very small* or close to *zero* reception rate even at very small distances (Fig. 3). The third major conceptual change is that there is a strong benefit of observing at least some percentage of links on-line. This is because some of the most effective links in terms of metrics of travel distance versus required number of messages are links that have a reception rate between 40-60%. In addition, we can observe from Figs. 3, 6(b) and 7 that it is perfectly possible to find high reception rate links that are stable and highly symmetrical that cover medium to long distances.

The complex and correlated nature of links implies that newly developed routing protocols should be simulated for much longer periods of time in order to ensure that overall they perform well. The existence of superior nodes in terms of both transmitters and receivers capabilities implies that fairness will become one of the major issues for any routing, multicast, and broadcasting approach, because all of these protocols have a tendency to disproportionately use a subset of nodes.

The statistically demonstrated space correlation will also greatly impact the development of routing protocols, as well as power management techniques. For example, since nodes are naturally clustered in subsets that efficiently communicate with each other and poorly with the rest of the network, it will be important that power management strategies, simultaneously turn down or up the majority of the nodes in one of such subsets. Furthermore, clustering techniques might be even more efficient than in networks modeled with the unit disk communication model.

In summary, we have developed a set of non-parametric statistical models for characterizing links in wireless sensor networks. The models are the basis for new generators of wireless networks to be used in simulations that are statistically similar to deployed networks. The insight gained while building these models has helped identifying future directions for developers of protocols and localized algorithms for wireless sensor networks.

IX. ACKNOWLEDGEMENTS

This work was made possible with support from The Center for Embedded Networked Sensing (CENS) under the NSF Cooperative Agreement CCR-0120778.

REFERENCES

- [1] A. Cerpa, N. Busek, and D. Estrin, "SCALE: A tool for simple connectivity assessment in lossy environments," CENS, UCLA, Tech. Rep. 0021, Sep 5 2003.
- [2] D. Ganesan, B. Krishnamachari, A. Woo, D. Culler, D. Estrin, and S. Wicker, "Complex behavior at scale: An experimental study of low-power wireless sensor networks," CENS, UCLA and IRL, UCB, Tech. Rep. 02-0013, February 2002.
- [3] Y. Zhao and R. Govindan, "Understanding packet delivery performance in dense wireless sensor networks," in *Proceedings of ACM Sensys 2003*. Los Angeles, CA, USA: ACM, Nov 5-7 2003, pp. 1-13.
- [4] Chipcon, "Cc1000 low power radio transceiver." [Online]. Available: <http://www.chipcon.com>
- [5] RFM, "Tr1000 low power radio system." [Online]. Available: <http://www.rfm.com>
- [6] W. Ye, J. Heidemann, and D. Estrin, "An energy-efficient MAC protocol for wireless sensor networks," in *Proceedings of the Twenty First Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM)*. New York, NY, USA: IEEE, June 2002, pp. 1567-1576. [Online]. Available: <http://www.isi.edu/johnh/PAPERS/Ye02a.html>
- [7] T. S. Rappaport, *Wireless Communication: Principles and Practice*. Prentice Hall, 2000.
- [8] A. Woo, T. Tong, and D. Culler, "Taming the underlying challenges of reliable multihop routing in sensor networks," in *Proceedings of ACM Sensys 2003*. Los Angeles, CA, USA: ACM, Nov 5-7 2003, pp. 14-27.
- [9] G. Zhou, T. He, S. Krishnamurthy, and J. A. Stankovic, "Impact of radio irregularity on wireless sensor networks," in *International Conference on Mobile Systems, Applications and Services*, 2004, pp. 125-138.
- [10] D. Aguayo, J. Bicket, S. Biswas, G. Judd, and R. Morris, "Link-level measurements from an 802.11b mesh network," in *Proceedings of ACM SIGCOMM 2004*. Portland, OR, USA: ACM, Aug 30-Sep 3 2004.
- [11] J. E. Gentle, W. Hardle, and Y. Mori, *Handbook of Computational Statistics, Concept and Methods*. Springer-Verlag, 2004.
- [12] J. W. Tukey, *Exploratory Data Analysis*. Addison-Wesley, 1977.
- [13] L. Girod, J. Elson, A. Cerpa, N. R. Thanos Stathopoulos, and D. Estrin, "Emstar: a software environment for developing and deploying wireless sensor networks," in *Proceedings of the 2004 USENIX Technical Conference*, Boston, Massachusetts, USA, June 27-July 2 2004.