

Data Segmentation for Region Detection in a Sensor Network

Anand Panangadan
Saban Research Institute
Childrens Hospital Los Angeles
4650 Sunset Blvd, Los Angeles, CA 90027, USA
Email: anandvp@usc.edu

Gaurav S. Sukhatme
Computer Science Department
University of Southern California
Los Angeles, CA 90089, USA
Email: gaurav@usc.edu

Abstract

We present an online method for automatically segmenting the data stream generated at a sensor node into distinct sub-sequences. The data at each node are assumed to be generated by a stochastic process. The data stream (sequence) sensed at each node is split into a new sub-sequence when the Kullback-Leibler distance between the probability distributions induced by the previously observed data and that induced by the most recently observed data exceeds a threshold. A distributed region growing algorithm is used to cluster spatially proximal sensor nodes that observe similar data streams. Thus, the space being monitored by the sensors can be divided into regions that observe similar activity. Continuous data segmentation and region growing enables changes in the pattern of activity to be observed as shifting spatial regions. We provide an analysis of the time taken for convergence and the power required for transmitting messages when the nodes are arranged in a regular grid and when they are randomly distributed. We present simulation results that show the performance of these algorithms on a shadow tracking application. The data used in these simulations was obtained from a sensor network equipped with light sensors and deployed outdoors.

I. INTRODUCTION

An application where sensor networks are increasingly being used is target tracking [1], [2]. In this application, a sensor network occupies the same physical space as the object being tracked. However, the presence of the object cannot be reliably detected using the sensory information gathered at one node. Instead, the impact of the object may be felt by multiple sensors that occupy the spatial region containing the object. Thus, it is of interest to map the extent of this region and report it to a central observer location. If this process is performed repeatedly, the observer will be able to monitor the changes in the location and extent of the region influenced by the object over time. Thus, region detection may be considered as a solution to the distributed spatial tracking problem.

The problem of boundary detection in sensor networks is also concerned with identifying the bounds of a region [3], [4]. However, a sensor node can determine if it is on a region boundary by querying only its neighboring nodes in the network. In the case of region detection, nodes that observe the same phenomenon are said to be in the same region only if these nodes are located in a contiguous spatial area. This is illustrated in Figure 1. Sensor nodes *A* and *B* observe the same phenomenon, as do all their neighbors, but they are in different regions. Thus, a distributed computation is required to compute the extent of a region.

In this work, we present a method for converting the stream of sensor measurements into a discrete set of phenomena. This enables the efficient comparison of the sensor behavior at different nodes. A distributed region growing algorithm is then used to spatially cluster sensor nodes that observe similar phenomena. Thus, the space being monitored by the sensors is divided into regions that observe similar activity. Changes in the spatial pattern of activity are reflected as changes in the spatial clustering of the nodes in that region.

The region detection problem is now described using a specific example and the steps of our proposed solution are described in the context of that example.

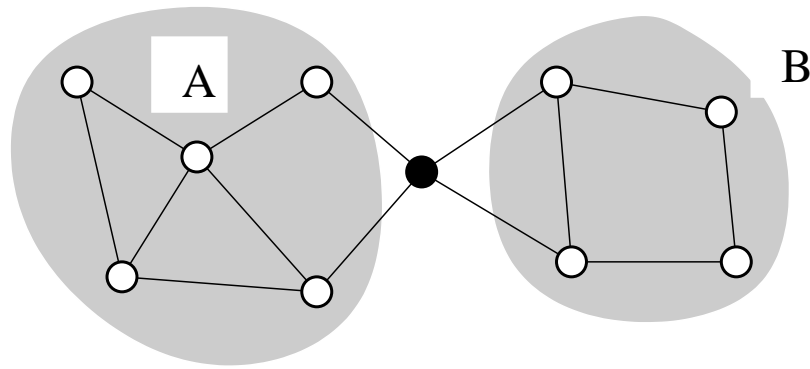


Fig. 1: Region detection requires distributed computation: Nodes *A* and *B* observe the same phenomenon along with all their neighbor nodes, but they belong to different regions.

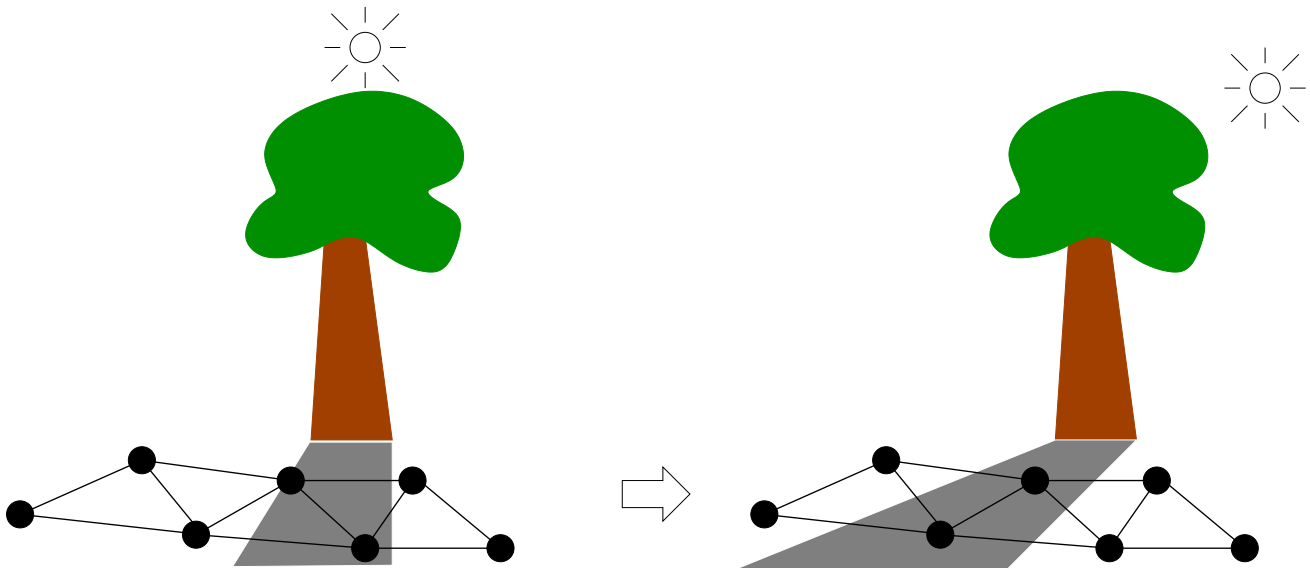


Fig. 2: A network of light sensors deployed to detect moving regions of shadow.

A. Motivating Example

Consider the task of tracking the movement of shadows on a forest floor with the aid of a light sensor network deployed on the forest floor. Each sensor is capable of determining whether its location falls under a shadow. The set of connected sensors that each detect a shadow together define a contiguous shadow region. The bounds of this shadow region should then be returned to a central query location in a succinct representation. As the shadow region moves during the day, the sensor network should be capable of detecting the change in the shape and location of the shadow region. The changing location and shape should be periodically available to the query location. This scenario is depicted in Figure 2.

To achieve this functionality, the following problems have to be solved.

- 1) Sensors that receive a stream of continuous-valued measurements should divide the data sequence into periods of distinct observed activity (such as light and shadow).
- 2) A distributed algorithm that enables the set of physically separate sensors to determine the bounds of the region that observes the same phenomenon. The algorithm should also be able to track the changing shape and location of this region.
- 3) A compact representation of the extents of the detected regions that can be transmitted to the query

location.

In this work, we first present a data segmentation algorithm to split the stream of sensor measurements into distinct sub-sequences. This algorithm considers measurements as arising from a stochastic process and uses a probability measure to detect a change in the distribution of the sensed data.

We then describe a distributed region growing algorithm that enables sensors observing the same phenomenon (as determined by the segmentation procedure) and belong to a connected component of the sensor network to converge on a common value. We also present a variation of the region growing algorithm that enables the sensor network to track changing region locations.

Rectangular bounding boxes are used to represent the extent of the regions detected by the region growing algorithm. We then describe how a query sent over the network returns the bounding box containing the locations of all the sensor nodes that comprise each region.

We provide an analysis of the time taken for convergence and the power required for transmitting messages when the nodes are arranged in a regular grid and when they are randomly distributed. We have implemented these algorithms on a network equipped with light sensors to track the progress of shadows during a day. We also present simulation results on larger-sized sensor networks.

II. DATA SEGMENTATION

A deployed sensor records data continuously. This stream of measurements may span more than one phenomenon. The task of segmentation is to divide this data sequence into segments each representing a distinct phenomenon. For instance, a light sensor detects high intensities when it is out of shadow but the intensity decreases when a shadow moves over it. These two periods of light and shadow should be detected as distinct phenomena. We use an information-theoretic metric to measure the distinctness of a data segment.

To apply this metric, we convert the sequence of continuous measurements to a probability distribution over a discrete space. The discretization is performed using vector quantization [5]. Vector quantization is a process of compressing a continuous input space into a discrete space. The discrete space is chosen so as to reflect the distribution of the data on the input space. Let X denote the set of input data measurements. We use the algorithm in Figure 3 to calculate N codewords that comprise the discrete space [6]: The algorithm initializes the codewords to N randomly picked input measurements. At every iteration, the current codewords are used to cluster the input data: each input measurement x is assigned to the cluster represented by the codeword that is closest to x . The codewords are recomputed to be the mean of every cluster. The iterations stop when the codewords converge to one set of values.

The algorithm to calculate the codewords can be run offline before the sensor network is deployed. The codewords are used to discretize the continuous measurements received by the sensors. Measurement x is discretized to that codeword that is closest to x .

We define an activity to be a set of measurements drawn from a fixed distribution. For instance, the probability distribution corresponding to daylight would have a high probability of high intensity measurements and low probability for the low intensity measurements. We assume that different phenomena give rise to distinct probability distributions of discretized measurements.

Each entire activity sequence is split into a number of consecutive subsequences such that these subsequences are distinct from each other in a probabilistic sense. We measure the difference between subsequences using an entropic measure known as the Kullback-Leibler (KL) distance [7]. Let p_1, p_2 be two probability distributions defined over the discrete space of measurements, \mathcal{X} (codewords). The KL distance between p_1, p_2 is defined as

$$KL(p_1, p_2) = \sum_{x \in \mathcal{X}} p_1(x) \log \left(\frac{p_1(x)}{p_2(x)} \right)$$

```

Algorithm calculate_codewords ( $N, X$ ) returns codewords
1: Initialize: codewords  $\leftarrow N$  random measurements from  $X$ 
2: repeat until codewords do not change
3:   for  $x \in X$ 
4:      $a[x] \leftarrow \operatorname{argmin}_{j=1,2,\dots,N} (|x - \text{codewords}_j|)$ 
5:   for  $j = 1, 2, \dots, N$ 
6:      $\text{codewords}_j = \sum_{x \in X, a[x]=j} x / \sum_{x \in X, a[x]=j} 1$ 
7: end repeat
8: return codewords

```

Fig. 3: Algorithm to calculate N codewords given a set of input measurements X .

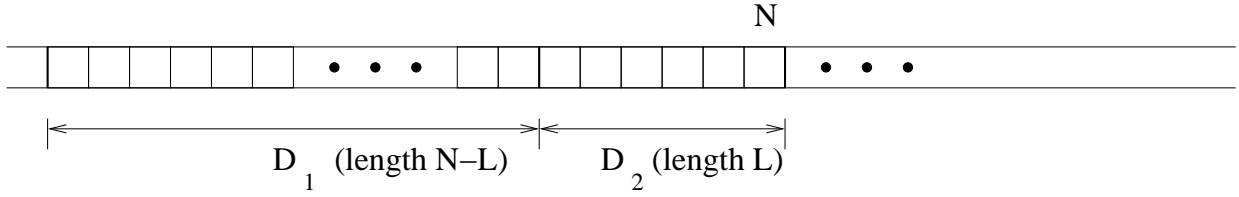


Fig. 4: The sliding window applied to the sequence of sensor measurements to determine if the sub-sequences D_1 and D_2 are generated by different probability distributions.

The KL distance between two distributions equals 0 only when the two distributions are equal. Let D denote a sequence of discrete measurements and $p(D)$ denote the Maximum Likelihood distribution obtained from D . The probability that a sequence D of length N is generated by a distribution p is approximately

$$e^{-N \cdot KL(p, p(D))}$$

Thus the probability that a sequence is not generated by distribution p drops exponentially with the length of the sequence.

The data sequence can be split into sub-sequences by recursively splitting the sequence at the point that gives rise to maximum divergence of the two component tracks. This method has been used for segmenting DNA sequences [8] and for splitting movement tracks of people for activity modeling [9]. The recursive method requires that the entire sequence be known before the splitting procedure can be applied. Here, we present an online segmenting procedure that can make the splitting decision as the data readings are obtained from the sensor.

Let $D = \{d_1, d_2, \dots, d_N\}$, $d_i \in \mathcal{X}$ be the sequence of observed measurements.

The online segmenting procedure utilizes a sliding window of length N .

D is split into two sub-sequences D_1 and D_2 :

$$D_1 = \{d_1, d_2, \dots, d_{N-L}\}, D_2 = \{d_{N-L+1}, d_{N-L+2}, \dots, d_N\}$$

This is illustrated in Figure 4. The two probability distributions, $p(D_1)$ and $p(D_2)$, are compared by calculating the KL distance between them. If $KL(p(D_1), p(D_2))$ is above a pre-defined threshold, then D is split at time $N-L$. Otherwise, the sliding window is moved forward when the new measurement d_{N+1} is obtained.

The online segmenting algorithm does not require that all the measurements comprising sub-sequence D_1 be stored. D_1 has to be represented only as its Maximum Likelihood distribution ($|\mathcal{X}|$ numbers) and

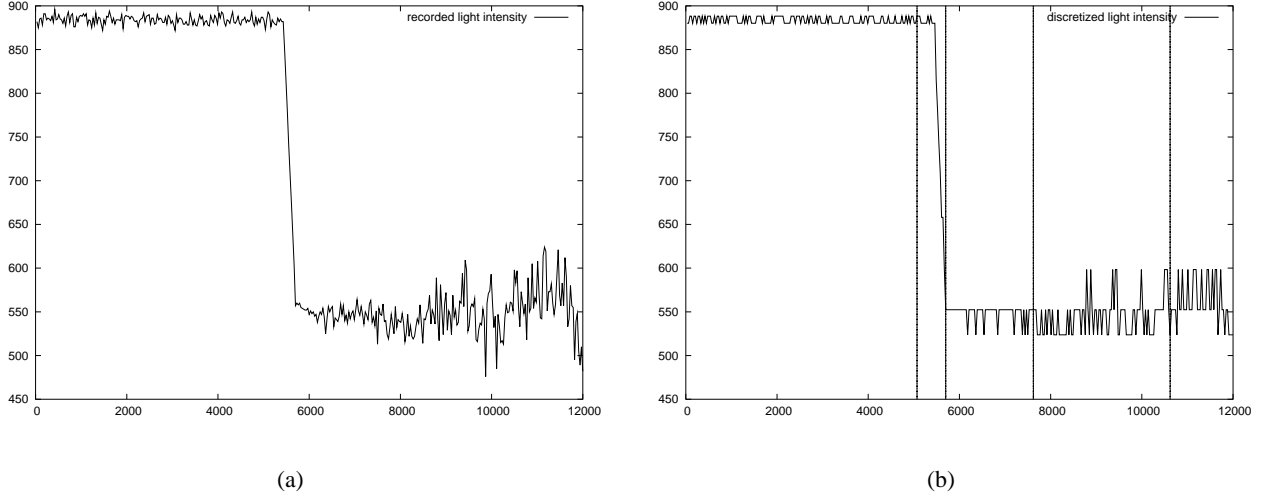


Fig. 5: (a) Light intensity recorded by a sensor node. x-axis shows time in seconds. The shadow crosses the node at $t = 5520s$. (b) Intensity values discretized into 10 values using vector quantization. The points where the sequence is split using the segmenting algorithm are marked by the vertical lines.

new measurements are used to update this distribution. Thus, this segmenting algorithm only requires that $|\mathcal{X}| + L$ numbers be stored at every step.

Figure 5 illustrates the segmentation procedure for a sequence of light intensity measurements (obtained from the setup described in Section IV). Figure 5(a) shows the raw continuous valued measurements plotted across time. These readings were used to generate 10 vector quantization codewords using the algorithm in Figure 3. The continuous intensity readings were then discretized by mapping each data point to the nearest codeword. The resulting sequence of discrete values is shown in Figure 5(b). We then apply the segmenting algorithm ($L = 10$) to determine those points where the sequence should be split into distinct sub-sequences. The five sub-sequences are indicated in Figure 5(b) by the vertical lines. The first sub-sequence corresponds to a period when the sensor was in bright light. The remaining sub-sequences cover periods where a shadow moved over the sensor and the intensity showed greater variation.

III. REGION DETECTION

Regions are contiguous divisions of space where a similar phenomenon is observed. For instance, the shadow cast by a structure will exhibit similar light intensities and hence the area covered by the shadow is a region. Region detection is the process by which a finite set of distributed point sensors approximate the bounds of a contiguous spatial region. Let s_1 and s_2 be two sensors. s_1 and s_2 are neighbors if they are within radio range, i.e, $d(s_1, s_2) < D$. Represent by $v(s_1)$ and $v(s_2)$ the phenomena recorded by s_1 and s_2 respectively. At any given instant of time, $v(s)$, is given by the probability distribution that currently describes the measurements being observed at s (Section II). Denote by $\text{nbd}(s)$ all neighbors of s that observe the same phenomena:

$$\text{nbd}(s) = \{s' | d(s, s') < D, v(s) = v(s')\}$$

s_1 and s_2 are defined to be in the same region if there exists a sequence of sensors $s_1, s'_1, s'_2, \dots, s'_n, s_2$ such that every pair of consecutive sensors in the sequence is a neighbor pair and $v(s_1) = v(s'_1) = v(s'_2) = \dots = v(s'_n) = v(s_2)$.

We now describe a region growing algorithm that enables a set of spatially distributed sensors to determine the extent of the region in which they are located. We assume that each sensor node s is

```

Algorithm Region_Growing ( $s$ ) returns  $R(s)$ 
1: Initialize:  $R(s) \leftarrow id(s)$ , Transmit  $R(s)$  to all  $s' \in nbd(s)$ 
2: while (true)
3:    $R \leftarrow$  set of received region identifiers
4:    $r^* = \min(R)$ 
5:   if ( $r^* < R(s)$ )
6:      $R(s) \leftarrow r^*$ 
7:     Transmit  $R(s)$  to all  $s' \in nbd(s)$ 
8:   end if
9: end while
10: return  $R(s)$ 

```

Fig. 6: Region Growing algorithm

aware of its geographic location, (x_s, y_s) . This location may be determined by distributed localization algorithms [10]. We also assume that each node is assigned a unique identifier, $id(s)$. If the sensors are aware of their geographic location, this location may be used as the identifier. The region growing algorithm assigns a *region identifier*, $R(s)$, to every sensor s . At every stage of the algorithm, sensors having the same region identifier are in the same region, i.e., if $R(s_1) = R(s_2)$ then s_1 and s_2 are within the same region. The region growing algorithm is presented in Figure 6.

The algorithm initializes each sensor to be in a separate region defined by its location (line 1). At every iteration, a sensor node sets its region identifier to the minimum of that of its neighbors and its current region identifier value (lines 3-5). On termination, all the nodes in a region have the same region identifier equal to the minimum of all the node identifiers in that region. If the region identifier is changed then the sensor transmits its new region identifier to all its neighbors (line 7). The transmissions end and the region identifiers stop getting updated when all the sensors within a region have the same region identifier.

The sensors in a region define a graph whose nodes represent the sensors and edges exist between nodes representing neighbor sensors. The node which has the minimum identifier is called the *root*. At every iteration of the region growing algorithm, denote the set of nodes whose region identifiers change to that of the root the *frontier*. Initially, the frontier consists of only the root node. Consider the case where all transmissions are reliable and iterations on sensors proceed in lock-step. At every iteration, the number of hops from the root to a frontier node increases by one. At the end of all transmissions, the frontier nodes correspond to those sensors that lie on the edge of the region. Thus, the number of iterations in the worst case equals the diameter of the graph.

An analysis of the running time of the algorithm and power consumed during transmissions when the nodes are arranged on a square grid is presented below. Let the number of nodes be $n \times n$ and let the graph be 4-connected, i.e., every non-edge node has four neighbors (Figure 7). The identifiers of the nodes are assigned randomly.

Running time: The diameter of this graph is $2n$ and thus the worst case running time of the algorithm occurs when the root is one of the corner nodes. The root's region identifier then has to reach the node at the corner opposite the root. Thus the time taken before all sensor nodes are assigned their region identifier value is $2n$.

Transmission power: Let the square grid have sides of unit dimensions. Then the distance between neighbor nodes is $r = 1/n$. Assume that the power required per transmission is proportional to the area of reception, i.e., kr^2 . At every region identifier update, a node s transmits $R(s)$ to all its neighbors. A sensor's region identifier does not change after it is set to that of the root. At the end of the i^{th} iteration,

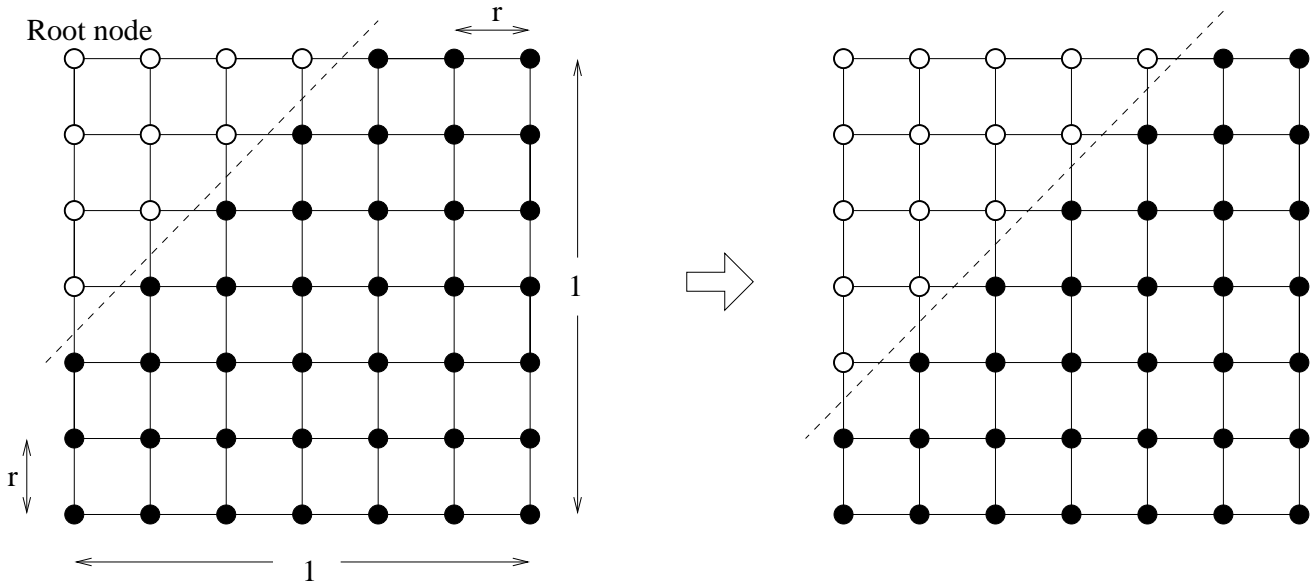


Fig. 7: Nodes on a unit square grid. The root for the region is a corner node. The unfilled circles indicate nodes that have already received the region identifier from the root node. This set of nodes increases by one hop in the next iteration. The dashed line indicates the frontier.

all sensors within i hops of the root have converged to their final value. Thus, the total transmission power is:

$$= kr^2 \left[\sum_{i=0}^{n-1} \left(n^2 - \sum_{j=0}^i j \right) + \sum_{i=1}^{n-1} \sum_{j=1}^i j \right] = kr^2 n^3 = kn$$

Thus, the total power consumed increases proportionally to the square root of the number of nodes deployed within the unit grid.

A. Dynamic Region Growing

The distributed region detection algorithm described above assumes that the regions being monitored do not change. This algorithm can be modified to handle the case when the region observing the phenomenon changes. The location and shape of the region may change either because the phenomenon has changed or because faults in some of the nodes can render previously contiguous regions discontinuous. In the algorithm in Figure 6, a node recalculates its region identifier only in response to region identifiers received from neighboring sensors.

To handle dynamic region detection, each node s stores the identifier of that neighbor that transmitted s 's region identifier. Denote this neighbor node, $\text{parent}(s)$. $\text{parent}(s) = s$ if node s is the root of a region. $\text{parent}(s)$ is set at line 6 in the region growing algorithm. At regular intervals, every sensor compares its region identifier with that of its parent. If the region identifiers do not match, or a response was not obtained from the parent node (node or communication fault), then node s resets its region identifier to its identifier: $R(s) \leftarrow id(s)$. This is equivalent to setting node s in a region by itself and re-initializing the region growing algorithm at node s .

B. Representing Regions Compactly

After the region identifiers have been assigned by the region growing algorithm, the extent of the regions can be calculated from a distributed query. Let the query to return the bounds of the region represented

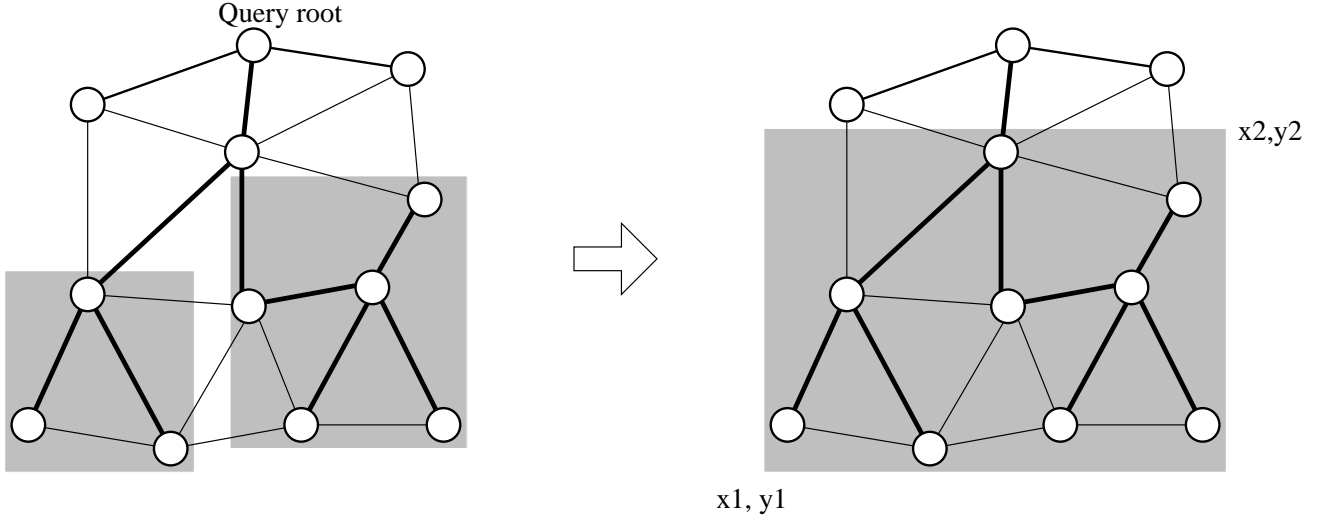


Fig. 8: Merging bounding boxes along the nodes of the query tree. Thin lines join neighboring sensor nodes. Thick lines indicate the edges of the query tree.

by region identifier r^* spread to all nodes in the sensor network along query tree T . Denote by $C_T(s)$ the set of child nodes of sensor node s on tree T .

We represent the extent of a region by a rectangular bounding box. Each sensor node s maintains the corners of the bounding box:

$$\text{bbox}(s) = (\text{bbox}_{x_1}(s), \text{bbox}_{y_1}(s), \text{bbox}_{x_2}(s), \text{bbox}_{y_2}(s))$$

of the region defined by the sensors rooted under s . The bounding boxes of a leaf node cover the physical location of the node.

$$\text{bbox}(s) = (x_s, y_s, x_s, y_s), C_T(s) = \phi$$

The bounding box of a non-leaf node s is obtained by merging the bounding boxes of those child sensors of s that have region identifier r^* . If $R(s) = r^*$,

$$\text{bbox}_{x_1}(s) = \min_{s' \in C_T(s), R(s')=r^*} (\text{bbox}_{x_1}(s'), x(s)); \text{bbox}_{y_1}(s) = \min_{s' \in C_T(s), R(s')=r^*} (\text{bbox}_{y_1}(s'), y(s))$$

$$\text{bbox}_{x_2}(s) = \max_{s' \in C_T(s), R(s')=r^*} (\text{bbox}_{x_2}(s'), x(s)); \text{bbox}_{y_2}(s) = \max_{s' \in C_T(s), R(s')=r^*} (\text{bbox}_{y_2}(s'), y(s))$$

If $R(s) \neq r^*$,

$$\text{bbox}_{x_1}(s) = \min_{s' \in C_T(s), R(s')=r^*} (\text{bbox}_{x_1}(s')); \text{bbox}_{y_1}(s) = \min_{s' \in C_T(s), R(s')=r^*} (\text{bbox}_{y_1}(s'))$$

$$\text{bbox}_{x_2}(s) = \max_{s' \in C_T(s), R(s')=r^*} (\text{bbox}_{x_2}(s')); \text{bbox}_{y_2}(s) = \max_{s' \in C_T(s), R(s')=r^*} (\text{bbox}_{y_2}(s'))$$

The bounding box at the root of the query tree represents the area defined by region identifier r^* . These steps are illustrated in Figure 8. Note that the region represented by a rectangular bounding box overestimates the actual physical area covered by the sensors. Also, the bounding boxes of two non-overlapping regions may overlap. Alternate representations of space for region detection are compared in [11].

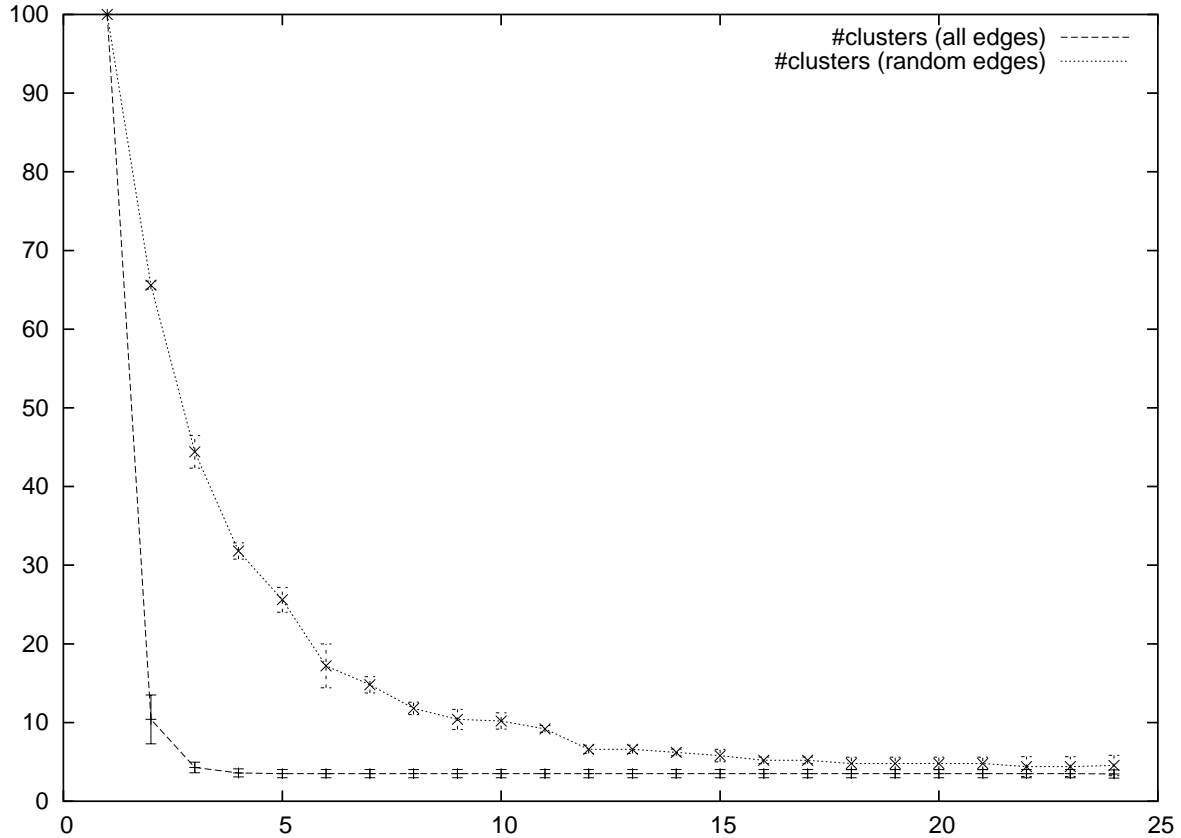


Fig. 9: Number of clusters detected after every iteration on 100 nodes randomly deployed on a square. The square was divided into three regions. The region growing algorithm was run with nodes considering all neighbors in every iteration, and with nodes only considering a random subset of neighbors in every iteration. Errorbars are at one standard deviation.

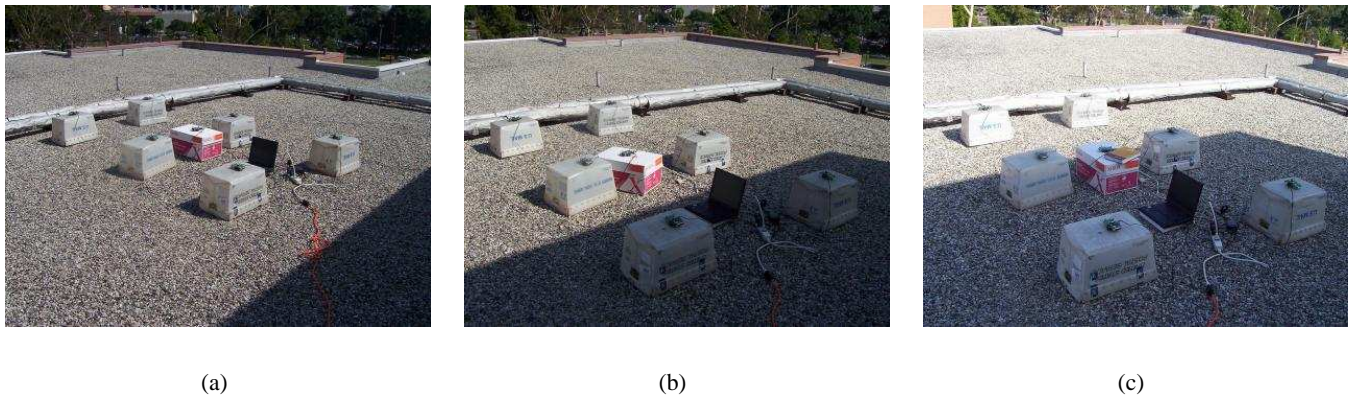


Fig. 10: Six motes deployed on the roof of a building. (a) All the motes are in direct sunlight. (b) A shadow has passed over two of the motes (toward the right). (c) Four of the motes are now in the shadow.

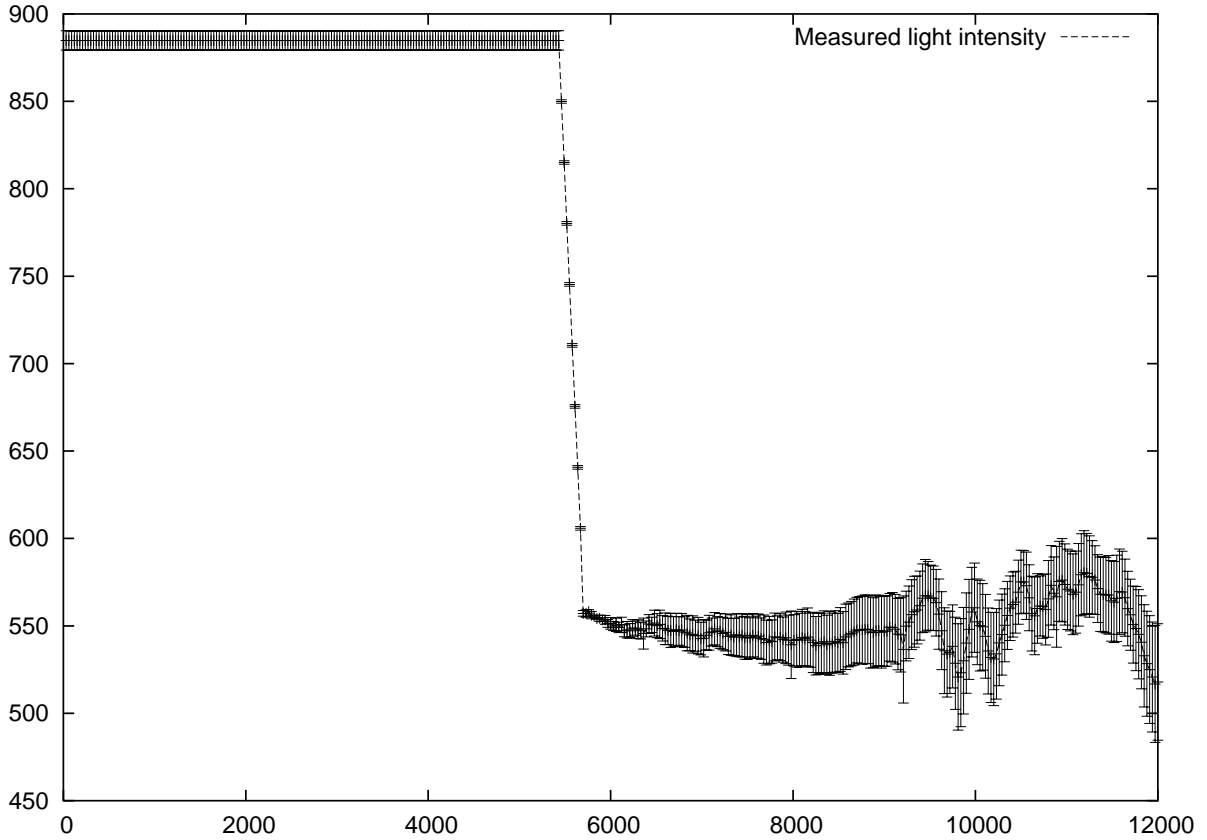


Fig. 11: The mean and standard deviation of the light intensities recorded by the six motes at every time interval. x-axis represents time in seconds. Errorbars are at one standard deviation.

IV. RESULTS

We deployed six mica2 motes equipped with light sensors on the roof of a building. The deployed area came into shadow after a certain time. The motes were arranged such that the shadow crossed them one after the other (Figure 10). The sensors recorded the light intensity every 30 seconds for approximately 3.5 hours. The data from one of the sensors is shown in Figure 5(a).

In order to run simulation trials for a larger number of motes based on this recorded data, we calculated the mean and variance at every time interval from the six data sequences (after adjusting for the initial offset to account for the difference in each mote’s position with respect to the shadow). The mean and standard deviation are shown in Figure 11. We were now able to generate simulated light intensity data for any physical deployment of sensor nodes. At any time instant, the light intensity reading of a simulated sensor would be drawn from a Gaussian distribution with mean and variance as calculated above. The sequence of generated data would be offset to account for the difference between the time that the shadow’s edge would cross the sensor and the actual shadow crossing time. Yu et al describe another method for extending a limited set of real sensor data to larger simulation models [12].

We deployed in simulation 100 nodes at randomly drawn locations on a square space of dimensions 100×100 . The area under light reported by the region growing algorithm as the virtual shadow moves across the deployed field is shown in Figure 12. The true area under light is also indicated. Note that the reported area under-estimates the true area because of edge effects. The match between the reported and true areas depends on the placement of the nodes near the boundaries of the region.

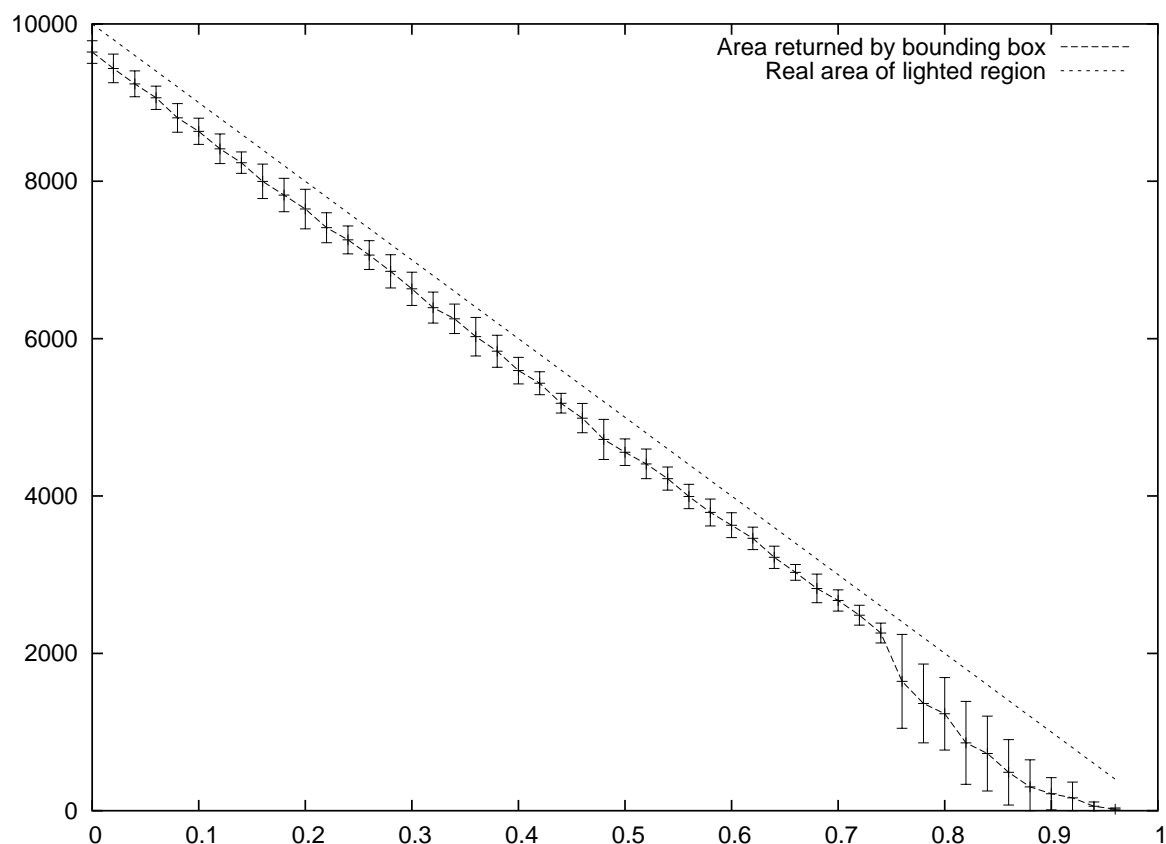


Fig. 12: The area under light reported by the region growing algorithm as a shadow moves across the deployed field. 100 nodes were randomly deployed in simulation on a square space of dimensions 100×100 . The intensity readings were obtained by randomly sampling from real sensor data. Errorbars are at one standard deviation. The true lighted area is also plotted.

V. CONCLUSIONS

We presented an online method for segmenting a continuous data stream generated at a sensor node into distinct sub-sequences. This method considers the sensor data to be generated by a stochastic process and uses a probabilistic measure to determine when a new phenomenon is observed. We described a distributed region growing algorithm to spatially cluster sensors that record similar data streams and thus divide the space being monitored by the sensors into regions observing a common phenomenon. We provided an analysis of the time taken for convergence and the power required for transmitting messages when the nodes are arranged in a regular grid and when they are randomly distributed. We presented simulation results using data obtained from a sensor network equipped with light sensors. In the simulation, randomly deployed sensors used this data to track shadows moving over the space occupied by the network.

REFERENCES

- [1] D. McErlean and S. Narayanan, "Distributed detection and tracking in sensor networks," in *Asilomar Conference on Signals, Systems, and Computers*, 2002.
- [2] S. Pattem, S. Poduri, and B. Krishnamachari, "Energy-quality tradeoffs for target tracking in wireless sensor networks," in *2nd Workshop on Information Processing in Sensor Networks (IPSN)*, 2003.
- [3] K. Chintalapudi and R. Govindan, "Localized edge detection in sensor fields," *Ad-hoc Networks Journal*, 2003.
- [4] R. Nowak and U. Mitra, "Boundary estimation in sensor networks: Theory and methods," in *2nd Workshop on Information Processing in Sensor Networks (IPSN)*, 2003.
- [5] A. Gersho and R. M. Gray, *Vector Quantization and Signal Compression*. Kluwer Academic Publishers, 1991.

- [6] Y. Linde, A. Buzo, and R. M. Gray, "An algorithm for vector quantizer design," *IEEE Transactions on Communications*, vol. COM-28, pp. 84–95, 1980.
- [7] J. Lin, "Divergence measures based on the Shannon entropy," *IEEE Transactions on Information Theory*, vol. 37, no. 1, pp. 145–151, 1991.
- [8] P. Bernaola-Galvan, R. Roman-Roldan, and J. Oliver, "Compositional segmentation and long-range fractal correlations in DNA sequences," *Physical Review E*, vol. 53, no. 5, pp. 5181–5189, 1996.
- [9] A. Panangadan, M. J. Matarić, and G. S. Sukhatme, "Detecting anomalous human interactions using laser range-finders," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE Press, Sep 2004, pp. 2136–2141.
- [10] D. Moore, J. Leonard, D. Rus, and S. Teller, "Robust distributed network localization with noisy range measurements," in *Proceedings of the Second ACM Conference on Embedded Networked Sensor Systems (SenSys)*, November 2004.
- [11] B. Greenstein, E. Kohler, D. Culler, and D. Estrin, "Distributed techniques for area computation in sensor networks," in *The First IEEE Workshop on Embedded Networked Sensors (EmNetS-I)*, 2004.
- [12] Y. Yu, D. Estrin, R. Govindan, and M. Rahimi, "Using more realistic data models to evaluate sensor network data processing algorithms," in *The First IEEE Workshop on Embedded Networked Sensors (EmNetS-I)*, 2004.