# A Balanced Parallel Clustering Protocol for Wireless Sensor Networks Using K-means Techniques

Liansheng Tan, Yanlin Gong and Gong Chen

*Abstract*—**For wireless sensor networks (WSNs), it is a challenging task how to schedule the energy resource to extend the network lifetime due to the fact that WSNs are usually powered by limited and non-rechargeable battery. A clustering scheme is helpful in reducing the energy consumption by aggregating data at intermediate sensor nodes. In this paper, we propose a balanced parallel K-means based clustering protocol; we term it BPK-means protocol. In this new protocol, we use K-means algorithm to cluster the sensor nodes, the cluster-heads are then selected in terms of two factors, they are a) the distance from node to cluster-center, and b) the residual energy. BPK-means only requires local communications: each tentative cluster-head only communicates with their topologically neighboring nodes and other tentative cluster-heads when achieving a distributed clustering scheme. The algorithm thus has the attractive feature of parallel computations. Moreover, BPK-means further balances the clusters to improve intra-cluster communication consumptions. We present the algorithm of this new protocol, analyze its computing properties, and validate the algorithm by simulations. Both theoretical analyses and simulation results demonstrate that BPK-means can achieve better load-balance and less energy consumptions when compared with LEACH. In addition, the BPK-means protocol is able to distribute energy dissipation evenly among the sensor nodes, which then prolong the system lifetime for the networks significantly.**

*Index Term*—**Wireless sensor network, K-means, clustering algorithm, network lifetime, throughput.**

## I. INTRODUCTION

Wireless sensor networks (WSNs) [1], [13] are usually formed by from a few to several thousands of sensor nodes and randomly scatter in a certain area. Because of the features such as random deployment and self-organizing, WSNs are deemed as the ideal candidates for a wide range of applications, such as target tracking and monitoring of critical infrastructures. But one of the most restrictive factors on the lifetime of WSNs is the limited energy resources of the deployed sensor nodes. Recent results [2], [16] show that hierarchical clustering routing algorithms have the better

adaptability and low energy consumption than flat routing algorithms for the large-scale sensor networks. In clustering algorithms, the networks are divided into clusters, and each cluster is formed by one cluster-head (CH) and several nodes. Therefore, the most important problem in clustering algorithm design is how to improve the structure of clusters and optimize the selection of CHs as these factors have the dominant impact on energy consumption and the lifetime of network.

This paper is organized as follows. Section 2 presents the related works. Section 3 describes the network model and our objectives. Section 4 proposes a balanced parallel K-means based clustering protocol. Section 5 analyzes the proposed clustering algorithm. Section 6 shows its effectiveness via simulations, and compares it with LEACH. Section 7 concludes the whole paper.

## II. RELATED WORKS

In order to increase energy efficiency and extend the lifetime of sensors, there have been substantial research interests in the area of clustering in WSN. We review some of the most relevant papers.

In [5] the authors propose a clustering protocol named LEACH. LEACH randomly selects a few sensor nodes as CHs and rotates this role to evenly distribute the energy load among the nodes in the network. The CHs are responsible for collecting data from other nodes and transmitting the aggregated data to the Base Station (BS). Although LEACH is able to increase the network lifetime, there are still a number of issues about this protocol. LEACH use the predetermined probability to decide the number of CHs in each round, there is the possibility that the number of elected CHs are overabundant or infrequent in some rounds. Moreover, it is not obvious that the CHs are uniformly distributed through the network; hence, some nodes will not have any CHs in their vicinity. It can't produce better structure of the clusters. Furthermore, due to the imbalance of clusters, some CHs must to take on the added work because there are many nodes in their clusters. Finally, the protocol assumes that all nodes begin with the same amount of energy capacity. It is not fit for non-uniform energy nodes.

An energy-aware [14], [17] variant of LEACH named LEACH-E is proposed in [6],in which the nodes with higher energy are more likely to become the CHs. But the imbalance for the structure of clusters and the distance between the

Liansheng Tan was with Department of Computer Science, Central China Normal University, Wuhan 430079, PR China. He is now with the Research School of Information Sciences and Engineering, The Australian National University, Canberra ACT 0200, Australia (e-mail: liansheng.tan@rsise.anu.edu.au).

Yanlin Gong and Gong Cheng are with Department of Computer Science, Central China Normal University, Wuhan 430079, PR China.

non-CH nodes and CHs may diminish the gain in energy consumption. LEACH-C [6] is a centralized clustering algorithm. It analytically determines the optimum number of CHs by taking into account the energy spent by all clusters. It needs global information to form better clusters.

PEGASIS [8], TEEN [10] and ATEEN [11] are based on LEACH protocol too. They improve the energy consumption by optimizing the data transmission pattern not by improving the cluster formation.

In [15], the authors study the theoretical aspects of the clustering problem and propose the clustering idea named balanced k-clustering. It theoretically analyzes that the energy optimization may be achieved by balancing the clusters and minimizing the transmission energy dissipation. It merely assumes that the balanced k-clustering problem can be solved by Voronoi diagram and min-cost flow and has no simulation results to prove it. It doesn't incorporate the dynamic and distributed nature.

From all the aforementioned protocols, we can find that the design of dynamic cluster formation directly affects the overall system energy dissipation. In this paper, we look at the improvement on energy consumption of the sensor node clustering problem. We present a balanced parallel K-means based clustering protocol and call it BPK-means. BPK-means protocol confirms the optimum number $K$ signifying the CHs in each round. The proposed protocol is based on K-means algorithm and optimizes the total spatial distance between sensor nodes and the cluster centers. This would reduce the energy dissipated by sensor nodes while making communication with the corresponding CHs. BPK-means balances each cluster in terms of number of sensor nodes. It would help in balancing the system load on each CH. Parallel computation and local communication would achieve a distributed clustering scheme and hence enhance scalability.

## III. PRELIMINARIES

We first describe the network model and then give our objectives.

### A. Network Model

Let us consider the scenario that there are N sensor nodes distributed uniformly in an M×M region. We make some assumptions about the sensor nodes and the underlying network model:

1) These sensor nodes and the BS are all stationary after deployment. This is typical for sensor network applications.

2) There is only one BS located far from the sensing field.

3) Sensors are heterogeneous and have the different initial energy. Compared with LEACH, there is no limit to initial energy, so it is more applicable.

4) Each sensor node knows own geographical information.

5) To avoid the delay time too much, we use one-hop communication. The CHs directly communicate with the

sensor nodes or BS.

### B. Channel Transmission Model

We compute the energy consumption using the First Order Radio Model [5]. The equations are used to calculate transmission costs and receiving costs for an $L$-bit message and a distance $d$ are shown below:

$$E_{TX}(L,d) = \begin{cases} L(E_{elec} + \varepsilon_{fs}d^2) & d < d_0 \\ L(E_{elec} + \varepsilon_{mp}d^4) & d \geq d_0 \end{cases} \quad (1)$$
$$E_{RX}(L) = LE_{elec}.$$

In which, the electronics energy, $E_{elec}$, depends on factors such as the digital coding, modulation, filtering, and spreading of the signal, whereas the amplifier energy, $_{fs}d^2$ or $_{mp}d^4$, depends on the distance to the receiver and the acceptable bit-error rate. $d_0$ is the distance threshold.

### C. The Objectives

The operation of BPK-means is broken into rounds. It is based on K-means algorithm and improve it. It uses parallel computing to get local cluster-center so as to reduce time complexity and achieve distributed clustering. Using BPK-means protocol, the sensor nodes in clusters are dense as possible. Moreover, it balances the clusters so that each cluster gets an average node number. In addition, the CHs are selected by two factors which are the distance from node to cluster-center and the residual energy.

These objectives can be stated as follows:

Given $N$ sensor nodes and $K$ cluster heads, we form sets $C_1, \cdots, C_K$ (clusters), $s_i$ is the sensor node in $C_j$, $|C_j|$ is the number of nodes in $C_j$, so $1 \leq i \leq N$, $1 \leq j \leq K$, such that:

1) All the clusters are balanced, i.e., the node numbers in the clusters satisfy

$$\frac{N}{K} - \delta \leq |C_j| \leq \frac{N}{K} + \delta, \quad 1 \leq j \leq K, \quad (2)$$

where is the unbalance factor. In this paper, our focus is on strictly balanced clusters and we assume $= 0$. With no loss of generality we assume here that $N$ is a multiple of $K$.

2) Our objective is to minimize the total spatial distance of the cluster structure.

$$D = \sum_{j=1}^{K} \sum_{s_i \in C_j} dist(loc(s_i), \overline{u_j}), \quad (3)$$

s.t. $\frac{N}{K} - \delta \leq |C_j| \leq \frac{N}{K} + \delta, \quad 1 \leq j \leq K,$

where $loc(s_i)$ and $\overline{u_j}$ are the locations of a sensor node and the cluster center in $C_j$. Function $dist$ is the Euclidean distance between a node and the corresponding center. Because communication is usually the main source of energy consumption in sensors and greatly depends on distance. Minimizing the total spatial distance within clusters can produce better cluster structure and reduce energy consumption.

3) The CHs are uniformly distributed over the sensor field while having relatively high residual energy and quite short

distance to the cluster centers.

## IV. BPK-MEANS CLUSTERING PROTOCOL

The BPK-means protocol consists of three components: calculating the optimum number of clusters by computing the energy consumption in each round approximately; dividing K balanced clusters to reduce intra-clusters communication expense by using K-means techniques; selecting the CHs by utilizing the cluster centers and the residual energy of nodes. The detailed description of the BPK-means is given in the following three subsections.

### A. The Optimum Number of Clusters

Suppose an $N$-node sensor network in the $M \times M$ square field with the BS located far away. The initial energy of the nodes randomly distributed in $[E_0, \alpha E_0]$ field. The component $E_0$ is the minimal initial energy and decides the maximal initial energy. We assume each node sends $L$-bit message to the cluster heads in each cluster during each round.

The initial total energy of the network defines as

$$E_{initial} = \sum_{i=1}^{N} E_i = \sum_{i=1}^{N} E_0 a_i = E_0 \sum_{i=1}^{N} a_i, \quad a_1 = 1. \quad (4)$$

If there are K clusters, there are on average $N/K$ nodes per cluster (one CH node and non-CH nodes). We define that:

$$Ave = N/K. \quad (5)$$

Then the total energy dissipated in a round [6] is

$$E_{round} = L(2NE_{elec} + NE_{DA} + K\varepsilon_{mp}d_{toBS}^4 + N\varepsilon_{fs}d_{toCH}^2), \quad (6)$$

In which, $E_{DA}$ is the energy dissipated in the CH for data gathering, $d_{toBS}$ is the average distance from the CHs to BS, and $d_{toCH}$ is the average distance from the nodes to CHs. Suppose the density of nodes is uniform throughout the cluster area, from [6] we have

$$d_{toCH} = \frac{M}{\sqrt{2\pi K}}, \quad (7)$$

where $M$ is the side of the given square field and $M^2$ denotes the area of this field. We can find the optimum number of clusters by setting the derivative of with respect to zero:

$$K_{opt} = \frac{\sqrt{N}}{\sqrt{2\pi}} \sqrt{\frac{\varepsilon_{fs}}{\varepsilon_{mp}}} \frac{M}{d_{toBS}^2}. \quad (8)$$

### B. Dividing and Balancing the Clusters

K-means algorithm [3], [4], [9], [12] is a clustering method, which applications range from unsupervised leaning of neural network, pattern recognitions, classification analysis, artificial intelligent, image processing, machine vision, etc. The basic principle of this algorithm is: given a set of $N$ D-dimension vectors without any prior knowledge about this set, the K-means clustering algorithm forms $K$ disjoint nonempty subsets $\{C_1, C_2, \cdots, C_K\}$ of vectors such that each vector $v_i$ ( $v_i \in C_j$, $1 \leq i \leq N$, $1 \leq j \leq K$ ) has the closest distance to cluster center $\overline{u_j}$, $1 \leq j \leq K$. The algorithm achieves this result by minimizing a square-error function $D$

such that the objective function

$$D = \sum_{j=1}^{K} \sum_{v_i \in C_j} dist(loc(v_i), \overline{u_j}).$$

is minimized.

ALGORITHM I
K-MEANS

| | K-means algorithm |
|---|---|
| 1 | Randomly select $K$ members of the set $N$ to form the initial value of $\overline{u_j}$, $1 \leq j \leq K$. |
| 2 | Compute distance $dist(v_i, \overline{u_j})$, $1 \leq i \leq N$, $1 \leq j \leq K$ of each vector such that $dist(v_i, \overline{u_j}) = \left| v_i - \overline{u_j} \right|$. |
| 3 | Choose vector members of the K subset according to their closest distance to $\overline{u_j}$, $1 \leq j \leq K$. |
| 4 | While $E$ is not stable: |
| 5 | Compute a new $\overline{u_j}$, $1 \leq j \leq K$ for each of the K subset. |
| 6 | Repeat the steps 2 and 3. |
| 7 | End |

The K-means algorithm has time complexity $O(K \times N \times I)$, where $K$ is the number of desired clusters, $N$ is the total number of training vectors and $I$ is the number of iterations.

K-means algorithm has been used as a clustering method in many application areas. It can assure the good effect for clustering. We expect that it can be used in WSN after improving. In order to achieve distributed clustering, each CH computes own cluster-center and communicates with other CHs, then modify the structure of cluster and finally form a stable cluster. In addition, the CHs balance the clusters to reduce the energy load of the CHs so as to prolong the network lifetime. Based on these thought, we propose a balanced parallel K-means clustering algorithm. It covers two important ideas which are the parallel K-means clustering [7] and the cluster-balanced step.

Now we are interested in developing parallel K-means by an advantage of wide availability and relatively lost-costs of distributed computing on a sensor network. The base idea is that the base station initiates the algorithm. It randomly generates initial K nodes as the tentative cluster-heads (TCHs). These nodes find all of its neighbors, and begin iteration one. Each TCH computes a center of its cluster, and then broadcast the cluster structure to other TCHs. After communicating with other TCHs, it computes and chooses the node member of the new cluster according to the closest distance to the new cluster center. The new cluster structures are formed by broadcast between the TCHs. Continuing the iteration process until square-error function $D$ is not stable or the number of iteration is more than the max number.

Then let us introduce the basic idea of the cluster-balanced step. Above we get the optimum number of clusters $K_{opt}$ if there are $N$ nodes. We define the average node number per cluster as $Ave$ in (5). The cluster-balanced step is added in each iteration process.

1) After clustering, the TCHs communicate with each

other about the node number.

2) We use $|C_1|$, $|C_2|$, …, $|C_K|$ respectively denote node number within the clusters, $|C_1| \geq |C_2| \ldots \geq |C_K|$, and the corresponding TCH is marked as TCH$_1$, TCH$_2$…TCH$_K$. If $|C_1| > Ave$, then TCH$_1$ initiates this course. If $|C_1| = |C_2| > Ave$, we can use some strategy decide only one TCH to do it. We suppose it is TCH$_1$.

3) TCH$_1$ computes $dist(s_i, \overline{u_j})$, where $s_i \in C_1$, $j \neq 1$ and $|C_j| < Ave$. This means to compute the distances between the nodes and each cluster center whose cluster's node number is less than $Ave$. TCH$_1$ get $s_i$ if its value is minimum, and mark its TCH as the corresponding TCH computed.

4) Repeating step2 and step3, and have $|C_1|$-$Ave$ adjusted, then $|C_1| = Ave$. TCH$_1$ balances its cluster.

5) TCH$_1$ send adjusted nodes' information to corresponding TCHs and notify TCH$_2$ the newest node numbers after adjusting.

6) Each TCH whose node number is more than Ave adjusts in turn. So it will achieve a balanced situation finally.

**Example 1:** Fig. 1 and Fig. 2 give an example to illustrate how the cluster-balanced step works. In this example, if we use K-means to cluster the sensor nodes, the sensor nodes $s_1$ to $s_3$ should belong to the $C_1$, the sensor nodes $s_4$ to $s_6$ should belong to $C_2$, and the sensor nodes $s_7$ to $s_{12}$ should belong to $C_3$. However, the nodes in $C_3$ are more than $Ave$ (=4). Thus this scenario motivates the cluster-balanced step. To have a balance among all the clusters, in our method we suggest that the sensor $s_7$ should be grouped into $C_2$ because its distance is nearest to $C_2$ and the sensor $s_{12}$ should be grouped into $C_1$ because its distance is the nearest to $C_1$.
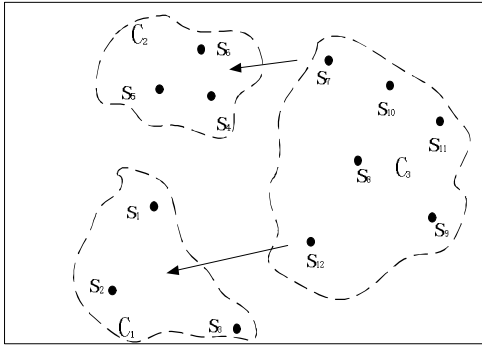

Fig.1 The clusters before balancing

### C. Selecting the CHs

After clustering and balancing, each TCH gets the cluster-center based on the information of sensor nodes within the cluster and then broadcasts it to these nodes. If we don't take into account the energy, we can choose the CH whose distance is nearest from itself to the corresponding cluster-center, but it isn't enough. Because the nodes closed to the cluster-center are dead earlier and the nodes with more energy far from the cluster-center aren't selected as the CHs.
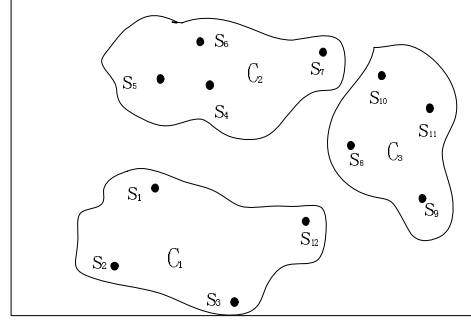

Fig.2 The clusters after balancing

Therefore, we can get the final balanced cluster structure by distributed computing and communicating.

ALGORITHM 2
BPK-MEANS

| | |
|---|---|
| The base station process: | |
| 1 | Randomly select $K$ TCHs as $K$ cluster centers |
| 2 | Send the average node number $Ave$ to each TCH. |
| The TCH$_j$ of cluster C$_j$ ($1 \leq j \leq K$) process: | |
| 1 | Broadcast a message to find its neighboring nodes. The Nodes decide to belong to which cluster according to the signal strength. The TCH$_j$ received join messages from its cluster members including the node ID and geographical information. |
| 2 | Implement the cluster-balanced step. |
| 3 | Compute $\overline{u_j}$ of $C_j$. |
| 4 | Broadcast $\overline{u_j}$ to every other TCH. |
| 5 | Compute distance $dist(s_i, \overline{u_j})$, $s_i \in C_j$. |
| 6 | Choose cluster member of the $K$ subsets according to their closest distance to $\overline{u_j}$. |
| 7 | Broadcast $K$ subset computed in step 6) to every other TCH. |
| 8 | Form the new clusters by collecting nodes that belong to $C_j$ that are sent from other TCHs in step 7. |
| 9 | While $D$ is not stable or loop number $\leq$ Maxloop |
| 10 | Repeat the steps 2 to 8 |
| 11 | End |

So the CHs are selected by two factors which are a) the distance from node to cluster-center, and b) the residual energy. The following function $T(d_{toCH}(s_i), \omega_i)$ can be used to decide node $s_i$ in its cluster to be the CH:

$$T(d_{toCH}(s_i), \omega_i) = d_{toCH}(s_i) \times exp(\frac{1}{\omega_i}), \quad \omega_i = \frac{E_i(r)}{\overline{E(r)}}, \quad (9)$$

In which, $r$ denotes the current round. $d_{toCH}(s_i)$ is the distance from node $s_i$ to its cluster-center. $\overline{E(r)}$ is the average residual energy given by

$$\overline{E(r)} = E_{total}(r) / N. \quad (10)$$

We compute the residual energy of the $r$-th round approximately for reducing the network traffic:

$$E_{total}(r) = E_{initial} - (r-1)E_{round} \quad (11)$$

Sensor node $s_i$ send the value computed from $T(d_{toCH}(s_i), \omega_i)$ to its TCH. Each TCH finds the corresponding node whose value is minimum and decide this node as the final CHs. The final CHs will notify each node about this, and then the communication begins.

The BPK-means protocol is observed to have the following properties:

**Theorem 1**: For the set of all sensors nodes $s=\{s_1 \cdots s_N\}$, the generated clusters $C_1$, ..., $C_K$ by the BPK-means algorithm have the following property：$|C_1| = |C_2| = \ldots = |C_K| = N/K =$ *Ave*. ( Again We assume that $N$ is a multiple of K ).

*Proof*: If $|C_1|$ is more than *Ave*, according to BPK-means Algorithm, it must implement the cluster-balanced step, so $|C_1|$ will equal to *Ave*. After each iteration process, $|C_1| = |C_2| = \ldots = |C_K| = Ave$, so when exit the iteration process, $|C_1| = |C_2| = \ldots = |C_K| = Ave$.

**Theorem 2**: Subject to the constraints $|C_1| = |C_2| = \ldots = |C_K| = Ave$, under the BPK-means algorithm the objective function

$$D = \sum_{j=1}^{K} \sum_{s_i \in C_j} dist(loc(s_i), \overline{u_j}),$$

is minimized approximately.

*Proof:* BPK-means is based on K-means. The objective of K-means clustering is to minimize the sum of the distances between all training vectors and their closest cluster centers. So without the constraint, D is minimized.

Now let us consider with the constraint. We check two cases. First, if $s_i$ is a membership in $C_j$ and the distance between $s_i$ and $\overline{u_j}$ is shorter than the distance between $s_i$ and other cluster centers, $D$ must be minimize. Second, because of the cluster-balanced step, $s_i$ becomes a member in $C_j$. So the distance between $s_i$ and $\overline{u_j}$ is shorter than the distance between $s_i$ and other cluster centers except for the cluster center which $s_i$ once belongs to. So with the constraint of the clusters being balanced, $D$ is near- minimized.

We demonstrate the above statement about the benefit in obtaining the suboptimal solution by the following numerical example.

**Example 2:** For the sensor network with 100 sensor nodes, we implement LEACH, K-means and BPK-means to cluster the 100 sensor nodes, and get the total spatial distance $D$ under these algorithms, respectively. Table 1 displays the network lifetime (in terms of the first node becoming dead time) and the resulted total spatial distance of the senor node structure under these three algorithms.

TABLE I
NUMERICAL EXAMPLE

|  | LEACH | K-means | BPK-means |
|---|---|---|---|
| D | 96627.738 | 34392.980 | 40080.534 |
| The first node dead time | 39 rounds | 98 rounds | 113 rounds |

From Table 1 we can find that K-means algorithm can get a minimum $D$. Because of the cluster-balanced step, BPK-means can get a bigger $D$ but the sensor nodes can survive a longer time. This implies that one can reach a certain tradeoff between the total spatial distance of sensor structure and the network lifetime. The suboptimal solution in BPK-means can achieve such tradeoff.

**Theorem 3**: Cluster heads are well-distributed over the sensor field while having relatively high residual energy and quite short distance to the cluster centers.

*Proof:* According our method for choosing the CHs, the CHs are closed to the cluster centers. Because the cluster centers are well-distributed over the sensor field by computing, the cluster heads do so.

Let us analyze (9), Function $T(d_{toCH}(s_i), \omega_i)$ is directly proportional to $d_{toCH}(s_i)$, when the residual energy of nodes are much at one, the probability that the nodes closed to the cluster centers are the CHs must be increased. On the other hand, when the distance between nodes and centers are much at one, the probability that the nodes hold higher residual energy are the CHs must be increased. Therefore, we can select some nodes with more energy and near the cluster center to be the CHs.

## VI. SIMULATION RESULTS

To evaluate the performance of BPK-means, we simulate BPK-means and LEACH using a random 100-node network with C++ and Matlab. The BS is located at (50, 150) in a 100 $\times$ 100 m² field. Supposing $d_{toBS} = 100\ m$, we compute $K_{opt} = 4$.

TABLE 2
PARAMETERS

| PARAMETER | VALUE | DESCRIPTION |
|---|---|---|
| $E_0 / J$ | 5 | The min-initial energy |
| $a$ | 2 | The factor deciding the max-initial energy |
| $L / bit$ | 4000 | Sending Data |
| $d_0 / m$ | 70 | Distance threshold |
| $E_{elec} /(nJ \cdot b^{-1})$ | 50 | Energy consumption per bit when sending and receiving |
| $\varepsilon_{fs} /(pJ \cdot m^{-2})$ | 10 | Energy consumption about the free space model amplifier |
| $\varepsilon_{mp} /(pJ \cdot m^{-4})$ | 0.0013 | Energy consumption about the multi-path model amplifier |
| $E_{DA} /(nJ \cdot b^{-1})$ | 5 | Energy consumption about data fusion |

Fig. 3 and Fig. 4 show the clustering structure for using LEACH and BPK-means. Comparing in Fig. 3 and Fig. 4, one finds that each cluster size is the same but the cluster heads locate more closely to the cluster centers by using BPK-means This gives us an intuition that it is more efficient to balance the load of network and to even distribute the nodes among clusters by using BPK-means algorithm. The benefits of using BPK-means algorithm are further demonstrated in Figs. 5 and 6, where we compare the network performance of network lifetime and throughput under the BPK-means protocol with that under LEACH and LEACH-E. Fig. 5 shows the total number of nodes that remain alive over the simulation time. While the first dead node remains alive for a more long time in BPK-means, this is because BPK-means takes account into the structure of clusters and the residual energy of nodes. Fig. 6 shows that BPK-means sends much more data in the simulation time than LEACH and LEACH-E.
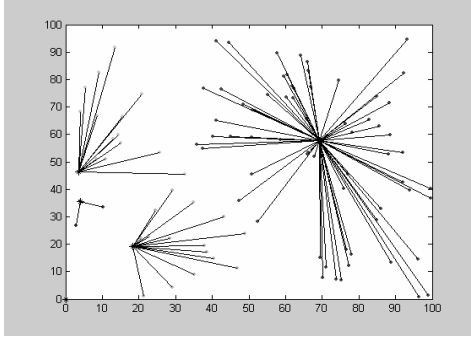
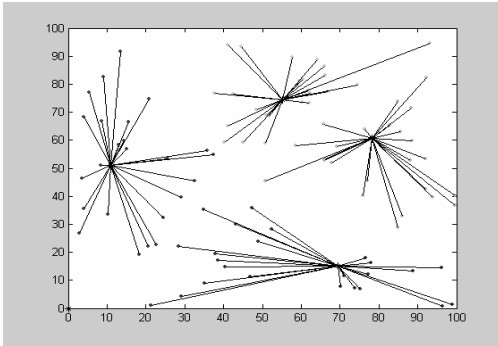Fig. 3 Dynamic clusters: the clustering structure for using LEACH



Fig. 4 Dynamic clusters: the clustering structure for using BPK-means
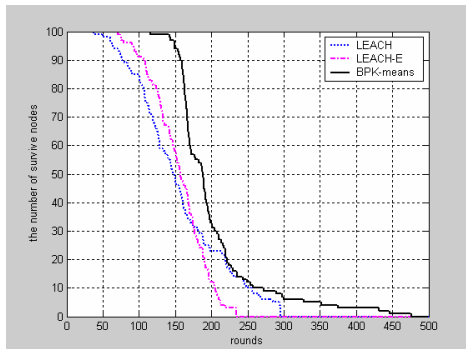

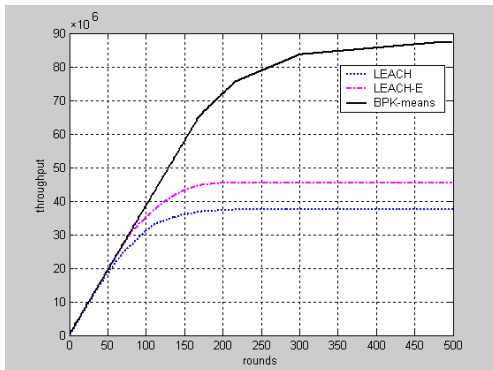
Fig. 5 System lifetime using LEACH, LEACH-E and BPK-means



Fig. 6 Throughput using LEACH, LEACH-E and BPK-means

## VII. CONCLUSIONS

In this paper, we propose the BPK-means clustering protocol for WSNs. BPK-means is based on a simple yet widely used method, namely K-means, which clusters objects based on attributes. BPK-means algorithm uses a distributed computing and cluster-balancedtechnique in achieving its clustering goal. Therefore, it can evenly distribute the energy load among all the sensor nodes to extend network lifetime. We provide the complete algorithm and give the theoretical analyses in demonstrating its application merit. The simulation results show that BPK-means significantly reduces the energy consumption, prolongs the lifetime of network and improve network throughput when compared with LEACH and LEACH-E.

### REFERENCES

[1] L. Akyildiz, W. L. Su, Y. Sankarasubramaniam and E. Cayirci, "A survey on sensor networks," *IEEE Commun. Mag.,* vol. 40, no. 8, pp. 102-114, Aug. 2002.
[2] J. N. Al-Karaki and A. E. Kamal, "A survey on routing protocol for wireless sensor networks," *IEEE Wireless Commun.,* vol. 11, no. 6, pp. 6-28, Dec. 2004.
[3] S. Datta, C. Giannella and H. Kargupta, "K-means clustering over a large, dynamic network," *Proc. 2006 SIAM Conf. Data Mining (SDM 06),* pp. 153-164, Apr. 2006.
[4] J. Ham and M. Kamber, "Data mining: concepts and techniques (2nd edition," *Morgan Kaufman Publishers,* pp. 1-6, 2006.
[5] W. Heinzelman, A. Chanrakasan, H. Balakrishnan, "Energy-efficient communication protocol for wireless micro-sensor networks," *Proc. of the 33rd Hawaii International Conf. on System Sciences*, vol. 2, pp. 1-10, Jan. 2000.
[6] W. Heinzelman, A. Chanrakasan, H. Balakrishnan, "An application-specific protocol architecture for wireless microsensor networks," *IEEE Trans. on Wireless Communications*, vol. 1, no. 4, pp. 660-670, Oct. 2002.
[7] S. Kantabutra and L. C. Alva, "Parallel K-means clustering algorithm on NOWs," *Technical Journal.* vol 1, no.5, pp. 243-249, 2000.
[8] S. Lindsey and C. Raghavendra, "PEGASIS: power-efficient gathering in sensor information systems," *Proc. IEEE Aerospace Conf.* vol. 3, pp. 9-16, 2002.
[9] S. P. Lloyd, "Least squares quantization in PCM," *IEEE Trans. Information Theory*, vol. 28, pp. 128-137, 1982.
[10] A. Manjeshwar and D. P. Agarwal, "TEEN: a routing protocol for enhanced efficiency in wireless sensor networks," *Parallel and Distributed Processing Symposium, Proceedings 15th International,* pp. 2009-2015, Apr. 2002.
[11] A. Manjeshwar and D. P. Agarwal, "APTEEN: a hybrid protocol for efficient routing and comprehensive information retrieval in wireless sensor networks," *Proc. of the International Parallel and Distributed Processing Symposium,* pp. 195-202, 2002.
[12] D. T. Pham, S. S. Dimov and C. D. Nguyen, "An incremental K-means algorithm," *Proc. Instn Mech. Engrs,* vol 218, pp. 783-795, Mar. 2004.
[13] G. J. Pottie and W. J. Kaiser, "Wireless Integrated Network Sensors," *Communications of the ACM*, vol. 43, no. 5, pp. 51−58, May 2000.
[14] L. Qing, Q. X. Zhu and M. W. Wang, "A distributed energy-efficient clustering algorithm for heterogeneous wireless sensor networks," *Journal of Software.* vol. 17, no. 3, pp. 481−489, Mar. 2006.
[15] G. Soheil, S. Ankur, J.Y.Xiao and M.Sarrafzadeh, "Optimal energy aware clustering in sensors networks," *Sensors,* vol. 2, no. 2, pp. 258-269, Jul. 2002.
[16] Y. Tang, M. T. Zhou and X. Zhang, "Overview of routing protocols in wireless sensor networks," *Journal of Software,* vol. 17, no. 3, pp. 410-421, Mar 2006.
[17] R. Wang, G. Z. Liu and Y. P. Shi, "Energy and distance efficient clustering algorithm for heterogeneous wireless sensor networks," *Journal of Wuhan University of Technology*, vol. 29, no. 4, pp. 110-113, Apr. 2007.