

Major Coefficients Recovery: a Compressed Data Gathering Scheme for Wireless Sensor Network

Liwen Xu, Yuexuan Wang, Yongcai Wang

Institute for Interdisciplinary Information Sciences

Tsinghua University, Beijing, 100084

kyoi.cn@gmail.com, {wangyuexuan,wangyc}@mail.tsinghua.edu.cn

Abstract—For large-scale sensor networks deployed for data gathering, energy efficiency is critical. Eliminating the data correlation is a promising technique for energy efficiency. Compressive Data Gathering (CDG) [8], which employs distributed coding to compress data correlation, is an important approach in this area. However, the CDG scheme uses a uniform pattern in data transmission, where all nodes transmit the same amount of data regardless of their hop distances to the sink, making it inefficient in saving transmission costs in 2-D networks. In this paper, the Major Coefficient Recovery (MCR) scheme is proposed, where the Discrete Cosine Transformation (DCT) is applied in a distributed fashion to the original sensed data. A non-uniform data transmission pattern is proposed by exploiting the energy concentration property of DCT and QR decomposition techniques so that sensors with larger hop-count can transmit fewer messages for network energy efficiency. The sink node recovers only the major coefficients of the DCT to reconstruct the original data accurately. MCR reduces the transmission overhead to $O(kn - k^2)$, an improvement by $O(\log n)$ over CDG in both 1-D and 2-D cases. The recovery performance of MCR is verified by extensive simulations.

I. INTRODUCTION

One of the major objectives of wireless sensor networks (WSN) is to collect information of physical phenomenon within a large-scale area. Energy efficiency and accurate observation of physical information are two goals of such data gathering networks. However, it is difficult to achieve both these goals at the same time, because accurate observation generally requires massive data gathering, posing great challenges to energy efficiency.

Considering the fact that many physical measurements are strongly correlated to the ones nearby [12], eliminating data correlation is a promising technique for energy efficient data collection. Such approaches [8][13] solve the wireless network data gathering problem by a distributed coding framework.

$$\mathbf{y} = \Phi \mathbf{d} = \Phi \Psi_{n \times n} \mathbf{x} \quad (1)$$

where $\mathbf{d} \in \mathbb{R}^{n \times 1}$ is a vector of nodes' original readings that can be compressed. $\Psi \in \mathbb{R}^{n \times n}$ is a selected basis, which maps \mathbf{d} to a sparse representation, i.e., in $\mathbf{d} = \Psi \mathbf{x}$. After compression, $\mathbf{x} \in \mathbb{R}^{n \times 1}$ is usually sparse. Φ is the *transmission pattern matrix* in data gathering networks, which determines how sensors transmit their data to the sink. \mathbf{y} is the observation at the sink. The objective is to accurately recover \mathbf{x} from \mathbf{y} .

For energy efficiency and accurate observation purposes, the problem in the above framework are: 1) How to select Ψ appropriately, so that the original observations can be compressively coded, and 2) How to optimize the design of Φ , so that the transmissions in a network can be reduced to prolong the network's lifetime.

Some work has been done to address the above problems. Since sensors work distributively, it is necessary for sensors to carry out compressive coding and transmission in a distributed fashion. Distributed Source Coding (DSC) was proposed in [13], where multiple correlated sensors compress their data distributively and send the compressed outputs to a central point for joint decoding. Based on Slepian-Wolf coding theory, [11] showed that distributed encoding can achieve the same efficiency as joint encoding can. This results ensure that the framework (1) can be carried out by distributed sensors.

Furthermore, for the selection of Ψ , Discrete Fourier Transform (DFT), Discrete Cosine Transform (DCT) and Wavelet Transform[3] is generally proposed as a transformation basis. For the optimization of Φ , Yoon et al. [14] proposed a clustered aggregation technique that first groups sensor nodes according to their measurements and transmits similar measurements per group only once in the process of data gathering. [4][5] proposed a scheme to encode relayed data for data gathering.

Most related to our work, a compressive data gathering (CDG) scheme was proposed by [8] for collecting data in tree structure networks. Sensors' observations are projected using random coefficients in Φ allowing each sensor to transmit the same number of messages, regardless of their hop distances to the sink. Such a uniform transmission pattern balances the loads of the sensors, but this results in a transmission overhead of $O(kn \log n)$: the same order as direct data collection takes without compressive sensing (non-CS) in 2-D networks. This is because when compared with non-CS, CDG assigns more work/heavier load for nodes further away from the sink and assigns less work/lighter load for nodes closer to the sink. In 2-D networks, sensors far away from the sink are more numerous than the sensors close to the sink, making CDG inefficient. [9] proposed hybrid compressive sensing (Hybrid-CS) to address such inefficiency, in which, a non-CS scheme is applied in the earlier stages of data collection (starting from the leaf nodes), and CDG is applied to nodes whose incoming traffic becomes greater than or equal to m , the number of rounds

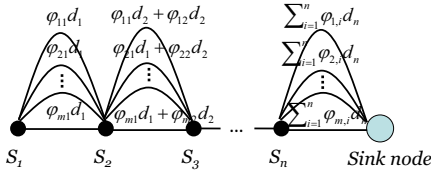


Fig. 1: Illustration of the relation between the transmission pattern and distributed source coding

of relay. But such a hybridized method needs to optimize the switching threshold and is not a unified framework for CS-based gathering. Furthermore, CDG and Hybrid-CS do not fully exploit the features of the sparse signal, which leaves space for further reduction of the number of transmissions.

In this paper, we exploit the energy concentration property of DCT-IV, and propose Major Coefficient Recovery (MCR) scheme to further reduce the number of transmissions at each sensor. In addition, we propose a QR decomposition-based algorithm to construct a non-uniform transmission pattern matrix for MCR, so that sensors far away from the sink can transmit fewer messages. MCR reduces the transmission overhead in 2-D network to $O(kn - k^2)$, which is an improvement by $O(\log n)$ over CDG. To improve recovery accuracy, Extended MCR is also proposed. The recovery performances of MCR are verified by extensive simulations.

The structure of this paper is organized as follows: Section II introduces the impact of transmission pattern matrices and design principles. Section III introduces the properties of DCT-IV, the proposed MCR scheme and the methods used in Extended MCR. Section IV presents some simulation experiment results. Section V concludes the paper with some discussion of future work.

II. TRANSMISSION PATTERN AND OVERHEAD

A. The relation between the transmission pattern matrix and transmission overhead

We define transmission overhead as the number of transmissions. A transmission pattern matrix determines how readings are coded and how they are transmitted in the network. In DSC, the original readings are encoded at each node, and then linearly combined with the data at each hop of the relay. Figure 2 illustrates the transmission pattern of a typical DSC in 1-D networks. s_1 has weight φ_{11} and transmits the product $\varphi_{11}d_1$ to node s_2 , where d_1 is the original measurement of s_1 . Then s_2 transmits the sum $\varphi_{11}d_1 + \varphi_{12}d_2$ to s_3 . At last, the sink node receives $\sum_{i=1}^n \varphi_{1i}d_i$ from s_n , i.e., a linear combination of the measurements of each node. Repeating this for m rounds the sink node receives a vector of combined measurements \mathbf{y} shown in (2).

$$\mathbf{y} = \Phi_{m \times n} \mathbf{d} = \Phi_{m \times n} \Psi_{n \times n} \mathbf{x} \quad (2)$$

The transmission pattern matrix determines the transmission overhead of a data gathering network. In one-dimensional cases:

- Each row of the transmission pattern matrix correlates to a single round of transmission. The number of trans-

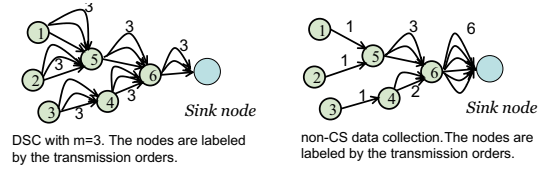


Fig. 2: Comparison of DSC and non-CS data collection. The numbers beside the links are the number of transmissions

missions in this round is equal to the number of entries that come after the first non-zero entry of this row. For example, if $\varphi_{i1} = \varphi_{i2} = \dots = \varphi_{ij} = 0$, then in the i -th round, nodes s_1, \dots, s_j don't need to transmit data.

The relation in two-dimensional cases is more complicated. However, under some specific settings, we can derive the relation similar to the one in 1-D cases. Consider the tree topology network as shown in Figure 2. Each node is labeled with a number, which implies the order of its relay. If it can be ensured that a node closer to the sink will not be labeled with a number smaller than the ones further from the sink, the relation between the transmission pattern matrix and the transmission overhead in 1-D networks still holds for 2-D tree structured networks.

It is crucial to design a suitable transmission pattern matrix so that the network is energy efficient. We can summarize two principles in designing the transmission pattern:

- **Fewer rows.** Fewer rows means fewer rounds of relay, which greatly reduces the total number of transmissions.
- **More leading zeros for each row.** Leading zeros in a row means fewer transmissions in this round of relay.

B. Transmission overhead of CDG analysis

CDG scheme proposed in [8] performs distributed coding, and employs compressive sensing technique to reconstruct original readings.

The transmission pattern Φ in CDG is an $m \times n$ random matrix with each entry obeying a normal distribution, satisfying the requirement of compressive sensing that Φ exhibits the restricted isometry property with high probability [1]. Ψ is chosen to be a DCT basis. Assuming the sparsity of \mathbf{x} is k , only if $m \geq ck \log n$, c is a positive constant, can \mathbf{x} be reconstructed with high probability according to the theory of compressive sensing [2]. So CDG only adopts the first design principle to optimize the transmission pattern matrix. The transmission overhead of CDG can be evaluated as follows:

Theorem 1. *The transmission overhead of CDG is at least of the order $O(kn \log n)$ if the topology of the network is a chain of n nodes or a tree with n nodes labeled by the transmission order as described above.*

Proof: The size of Φ is $m \times n$ which shows the number of transmissions is $O(mn)$. Since $m \geq ck \log n$, $mn \sim O(kn \log n)$. ■

Because the transmission overhead of non-CS data collection is $O(n^2)$ in 1-D network and is $O(n \log n)$ in 2-D networks[7], CDG reduces the transmission overhead in 1-D

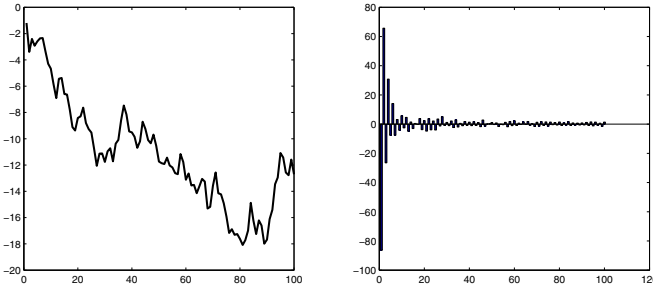


Fig. 3: A Markov-sequence-signal and its DCT

networks efficiently, but has the same order of overhead as the non-CS collection in the 2-D networks.

III. MAJOR COEFFICIENTS RECOVERY

When we define MCR as an extension of the distributed coding framework (1), then the first several *major coefficients* of \mathbf{x} can be recovered:

$$\mathbf{y} = \Phi_{k \times n} \mathbf{d} = \Phi_{k \times n} \Psi_{n \times n} \mathbf{x} = \mathbf{A}_{k \times n} \mathbf{x} \quad (3)$$

where k is the number of major coefficients that MCR tries to recover. Note that the size of transmission pattern matrix is $k \times n$, where k rounds of relay are carried out. Ψ in MCR is a basis of Discrete Cosine Transform type-IV (DCT-IV), whose properties of decorrelation and concentration allow MCR to achieve an accurate reconstruction by recovering only very few first coefficients of \mathbf{x} .

A. Energy Concentration Property of DCT-IV

There are several different kinds of DCT definitions, and DCT-IV is considered in this paper for its orthogonality, energy concentration and simplicity of definition. The transform matrix of a 1-D DCT-IV (of length n) is defined as

$$\Psi = \{\psi_{ij}\}, 1 \leq i, j \leq n$$

$$\psi_{ij} = \sqrt{\frac{2}{n}} \cos \left[\frac{\pi}{n} \left(i - 1 + \frac{1}{2} \right) \left(j - 1 + \frac{1}{2} \right) \right]$$

From the symmetry and orthogonality of Ψ , we know

$$\Psi = \Psi^T = \Psi^{-1} \quad (4)$$

DCT exhibits two important properties:

- 1) *Decorrelation* which can transform the correlated signal into unrelated coefficients that are usually sparse.
- 2) *Energy concentration*. The DCT coefficients are the combinations of cosine values of different frequencies. [10][6] illustrated that DCT coefficients tend to obey a Laplacian distribution in most practical cases. As k , i.e., the frequency, gets higher, the scale parameter b of a Laplacian distribution gets smaller, which means the corresponding DCT coefficient has a greater probability of taking on a very small value, even zero.

Figure 3 is the DCT representation of a typical signal. Nearly 95% of the total energy is concentrated within the first 10% of coefficients.

B. Transmission Pattern Matrix Design in MCR

To address the energy concentration property of DCT, we observe that the k -lowest frequency coefficients in \mathbf{x} dominate the expression of $\Psi_{n \times n} \mathbf{x}$. Therefore, a Major Coefficient Recovery scheme is proposed to recover \mathbf{d} using only the first k , i.e., the major coefficients, of \mathbf{x} , which substantially reduces the transmission cost.

The $k \times n$ matrix \mathbf{A} can be divided into two submatrices of size $k \times k$ and $k \times (n - k)$ denoted by $\mathbf{A}_k, \mathbf{A}_r$ in (5).

$$\mathbf{y} = (\mathbf{A}_k \ \mathbf{A}_r) \mathbf{x} \quad (5)$$

If \mathbf{A}_r is further made to be $\mathbf{0}$, (5) becomes

$$\mathbf{y} = (\mathbf{A}_k \ \mathbf{0}) \mathbf{x} \quad (6)$$

On the condition that \mathbf{A}_k is full rank, (6) turns out to be a linear system of equations whose solution is exactly the first k coefficients of \mathbf{x} .

The only requirement to ensure is that matrix \mathbf{A} consists of a $k \times k$ full-rank matrix and a zero matrix. Because Ψ is fixed as long as the number of nodes is fixed, we can only manipulate the transmission pattern for Φ in order to ensure that \mathbf{A} takes on the proper form. We propose an algorithm which can guarantee this.

Algorithm 1 Transmission Pattern Matrix Design in MCR

Input: n , the number of nodes in a line

Output: Φ , the transmission pattern

Step 1. Calculate $\Psi = (\psi_{ij})$, the DCT-IV basis of length n , where $\psi_{ij} = \sqrt{\frac{2}{n}} \cos \left[\frac{\pi}{n} \left(i - 1 + \frac{1}{2} \right) \left(j - 1 + \frac{1}{2} \right) \right]$

Step 2. Let Ψ_u be a $k \times n$ submatrix consisting of the first k rows of Ψ

Step 3. Do QR decomposition to Ψ_u , therefore $\Psi_u = \mathbf{Q}_{k \times k} \mathbf{R}_{k \times n}$, where \mathbf{Q} is an orthogonal matrix and \mathbf{R} is an upper triangular matrix

Step 4. Let $\Phi = \mathbf{R}$

Claim 1. Φ generated by Algorithm 1 guarantees the successful recovery of the k major coefficients of \mathbf{x} .

Proof: From the algorithm we know,

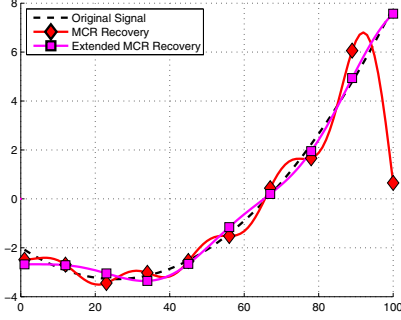
$$\begin{aligned} \because \Psi_u &= \mathbf{QR} = \mathbf{Q}\Phi \\ \therefore \Phi &= \mathbf{Q}^{-1} \Psi_u \\ &= (\mathbf{Q}^{-1} \ \mathbf{0}) \begin{pmatrix} \Psi_u \\ \Psi_l \end{pmatrix} \end{aligned}$$

where Ψ_l is an $(n - k) \times n$ submatrix consisting of the last $n - k$ rows of Ψ . So we have,

$$\begin{aligned} \Phi &= (\mathbf{Q}^{-1} \ \mathbf{0}) \Psi \\ \Phi \Psi^{-1} &= (\mathbf{Q}^{-1} \ \mathbf{0}) \end{aligned}$$

Since Ψ is symmetric and orthogonal as was presented in (4), we have

$$(\mathbf{Q}^{-1} \ \mathbf{0}) = \Phi \Psi = \mathbf{A}$$


 Fig. 4: Recovery of \mathbf{d} in MCR and extended MCR

\mathbf{A} is defined in (5), so we can see that Φ guarantees that \mathbf{A} consists of a $k \times k$ full-rank matrix and a zero matrix. ■

Put another way, the sink is to solve the linear system of equations

$$(\mathbf{Q}^{-1} \mathbf{0})\mathbf{x} = \mathbf{y} \quad (7)$$

C. Extended MCR

MCR recovers the major coefficients in DCT while ignoring the components of high-frequency. Though the value of high-frequency coefficients are usually very small, their loss may lead to divergence of recovery in the latter part of \mathbf{d} as is shown in Figure 4.

Since the former part of \mathbf{d} can be recovered accurately, we would be able to achieve a good global recovery accuracy if only “former” part could be “extended”. This can be achieved by inserting *virtual nodes* between the last node (the node next to the sink) and the sink node. Accordingly, *virtual measurements* can be generated and subsequently collected by the sink node. In the recovery process, the sink node can then simply throw out the virtual parts and keep the real measurements, which will improve the recover accuracy of MCR.

In practice, the simplest way to generate these virtual measurements is to repeat the last node’s reading, one for each virtual node we require. Figure 4 gives an example of this method which shows the increase in recovery accuracy.

D. Transmission Overhead of MCR

The transmission pattern matrix of MCR is a $k \times n$ upper triangular matrix.

$$\begin{pmatrix} \varphi_{11} & \varphi_{12} & \cdots & \varphi_{1k-1} & \varphi_{1k} & \cdots & \varphi_{1n} \\ 0 & \varphi_{22} & \cdots & \varphi_{2k-1} & \varphi_{2k} & \cdots & \varphi_{2n} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \varphi_{k-1k-1} & \varphi_{k-1k} & \cdots & \varphi_{k-1n} \\ 0 & 0 & \cdots & 0 & \varphi_{kk} & \cdots & \varphi_{kn} \end{pmatrix}$$

Theorem 1. *The number of transmissions is at least of the order $O(kn - k^2)$ if the topology of the network is a chain of n nodes or a tree with n nodes labeled the transmission order as described in Section II-A.*

Proof: Φ , the transmission pattern of MCR, is a $k \times n$ upper triangular matrix. If $i < j$, we have $\varphi_{ij} = 0$. So in the j -th round of data transmission, the nodes $\{1, 2, \dots, j-1\}$ do not

TABLE I: MSE and transmission overhead comparison

EXP #	MSE			Transmission Overhead	
	CDG	MCR	ExMCR	CDG	MCR/ExMCR
1	0.0058	0.0029	0.0049	5000	955
2	0.0159	0.0116	0.0129	5000	955
3	0.0168	0.0080	0.0051	5000	955
4	0.7783	0.6177	0.0291	20000	1955

need to transmit data. Therefore, the number of transmissions is the number of non-zero entries in Φ , i.e. $kn - \frac{k(k-1)}{2} \sim O(kn - k^2)$. ■

Compared with CDG, MCR is much more energy efficient. In both 1-D networks and 2-D tree structured networks, MCR’s transmission overhead is of the order $O(kn - k^2)$, an $O(\log n)$ improvement over CDG.

IV. EXPERIMENTS AND ANALYSIS

Our experiments are designed to compare the recovery accuracy and transmission overhead of MCR and the extended MCR (ExMCR for short) with CDG. The recovery accuracy is measured using Mean Square Error (MSE). According to the original signal types, the experiments are separated into 3 sets: polynomial signal, polynomial signal with white noise and Markov Chains. Table I shows the results of these experiments.

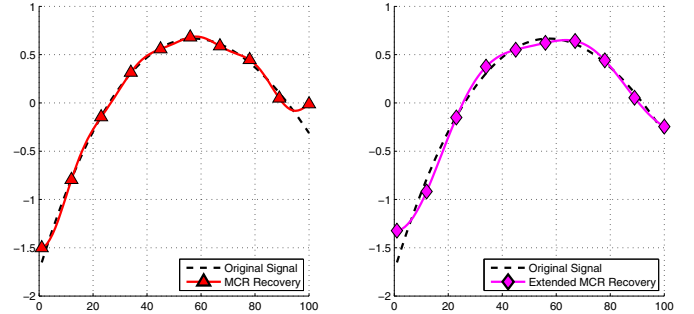
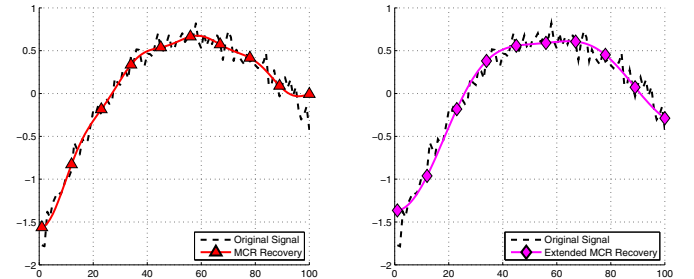


Fig. 5: EXP #1: Polynomial-signal


 Fig. 6: EXP #2: Polynomial-signal with noise $n \sim \mathcal{N}(0, 0.1)$

In the first 3 experiments, CDG, MCR and ExMCR all recover 10 out of 100 DCT coefficients. In the polynomial cases, CDG, MCR and ExMCR provide very accurate recoveries while MCR and ExMCR perform slightly better than CDG. In

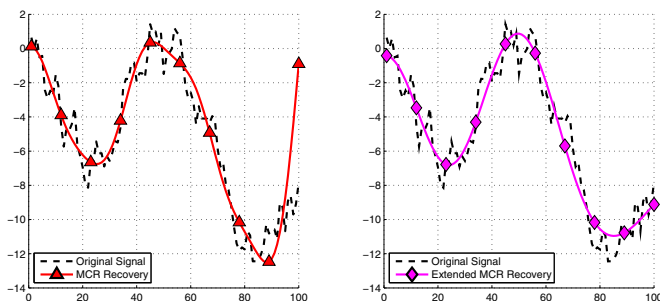


Fig. 7: EXP #3: Markov-sequence-signal

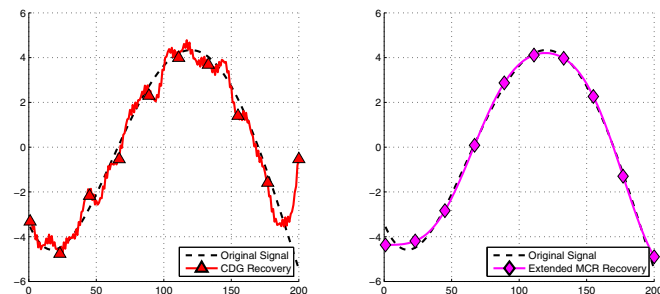


Fig. 8: EXP #4: Polynomial-signal recovered at different accuracy settings

experiment #3, the original signals are generated by Markov Chains. The experiment is repeated 50 times and the average values of MSE are presented. We can see MCR and ExMCR also perform better in this experiment. In experiment #4, a signal of 200 samples is generated by a polynomial function, and CDG tries to recover 15 out of 200 coefficients and MCR/ExMCR tries to recover 10. As shown in Figure 8, CDG recovery diverges from the original signal throughout, MCR begins to diverge in the latter part and ExMCR recovers better than both. The result of MSE also supports this observation.

We summarize the results of our experiments and highlight two points:

- **More accurate.** MCR and ExMCR provide accurate recoveries as good as or better than CDG in most cases at the same recovery setting.
- **Less overhead.** In CDG, 50 rounds of transmission are carried out in order to meet the requirement of compressive sensing, while in MCR, there are only 10 rounds of transmission, i.e. equal to the number of DCT coefficients we choose to recover. The number of transmissions in CDG is 5000, in MCR is 955, i.e. 19.1% of the number of transmissions in CDG. This overhead ratio is 9.78% in experiment #4.

V. CONCLUSION

In this paper, we propose the Major Coefficients Recovery, or MCR, scheme for data gathering in wireless sensor network. MCR only recovers the first k coefficients that are considered as major because they are likely to be larger than the rest with high probability according to the energy concentration

property of DCT. We further propose Extended MCR in order to avoid the divergence in data recovery that MCR suffers. Compared with CDG, the number of transmissions is greatly reduced in MCR. The transmission overhead is $O(kn - k^2)$ in both 1-D networks and 2-D tree structured networks, an $O(\log n)$ improvement over CDG. In the future, this work can be extended by distributing the DCT transformation and exploring the node numbering strategy in a greater number of 2-D topologies.

ACKNOWLEDGMENT

This work is supported in part by the National Basic Research Program of China Grant 2007CB807900, 2007CB807901, the National Natural Science Foundation of China Grant 61073174, 61033001, 61061130540, and the Hi-Tech research and Development Program of China Grant 2006AA10Z216.

REFERENCES

- [1] E. Candès, T. Tao. *Near-optimal signal recovery from random projections and universal encoding strategies*. IEEE Trans. Inform. Theory 52:489-509, 2006.
- [2] E. Candès. *Compressive sampling*. In Proc. of ICM, 2006.
- [3] A. Ciancio, S. Patten, A. Ortega, and B. Krishnamachari. *Energy-efficient data representation and routing for wireless sensor networks based on a distributed wavelet compression algorithm*. In Proc. of IPSN, pages 309-316, 2006.
- [4] R. Cristescu, B. Beferull-Lozano, and M. Vetterli. *On network correlated data gathering*. In Proc. of IEEE Infocom, volume 4, pages 2571-2582, Mar. 2004.
- [5] R. Cristescu, B. Beferull-Lozano, M. Vetterli, and R. Wattenhofer. *Network correlated data gathering with explicit communication: Np-completeness and algorithms*. IEEE/ACM Trans. on Networking, 14(1):41-54, Feb. 2006.
- [6] E. Y. Lam and J. W. Goodman, *A Mathematical Analysis of the DCT Coefficient Distributions for Images*, IEEE Trans. Image Process., vol. 9, no. 10, pp. 1661-1666, Oct. 2000.
- [7] X. Li, Y. Wang, and Y. Wang, *Complexity of Data Collection, Aggregation, and Selection for Wireless Sensor Networks*, IEEE Tran. on Computers, vol. 60 no. 3, pp. 386-399, 2010
- [8] C Luo, F Wu, J Sun, *Compressive Data Gathering for Large-scale Wireless Sensor Networks*, in Proc. ACM Mobicom'09, pp. 145-156, Sep. 2009.
- [9] J. Luo, L. Xiang, and C. Rosenberg, *Does Compressed Sensing Improve the Throughput of Wireless Sensor Networks?*, ICC'09, pp.1-6, 2009
- [10] R. C. Reininger and J. D. Gibson, *Distributions of the Two-dimensional DCT Coefficients for Images*, IEEE Trans. Commun., vol. COM-31, no. 6, pp. 835-839, Jun. 1983.
- [11] D. Slepian and J. K. Wolf. *Noiseless encoding of correlated information sources*. 19:471-480, Jul. 1973.
- [12] M. C. Vuran, Ö. B. Akan, and I. F. Akyildiz, *Spatio-temporal Correlation: Theory and Applications for Wireless Sensor Networks*, Computer Networks Journal (Elsevier), vol. 45, no. 3, June 2004.
- [13] Z. Xiong, A. Liveris, and S. Cheng, *Distributed source coding for sensor networks*, IEEE Signal Processing Mag., vol. 21, pp. 80C94, Sept. 2004.
- [14] S. Yoon and C. Shahabi. *The Clustered Aggregation (CAG) Technique Leveraging Spatial and Temporal Correlations in Wireless Sensor Networks*. ACM Trans. on Sensor Networks, 3(1), Mar. 2007.