

# Wireless Sensor Network

**Chief Editor : Kosai Raoof**



# Journal Editorial Board

ISSN 1945-3078 (Print) ISSN 1945-3086 (Online)

<http://www.scirp.org/journal/wsn/>

---

## Editor-in-Chief

**Dr. Kosai Raoof** University of Joseph Fourier, Grenoble, France

## Editorial Board (According to Alphabet)

<b>Prof. Dharma P. Agrawal</b>	University of Cincinnati, USA
<b>Prof. Ji Chen</b>	University of Houston, USA
<b>Dr. Yuanzhu Peter Chen</b>	Memorial University of Newfoundland, Canada
<b>Prof. Jong-wha Chong</b>	Hanyang University, Korea (South)
<b>Prof. Laurie Cuthbert</b>	University of London at Queen Mary, UK
<b>Prof. Thorsten Herfet</b>	Saarland University, Germany
<b>Dr. Li Huang</b>	Stiching IMEC Netherlands, Netherlands
<b>Dr. Yi Huang</b>	University of Liverpool, UK
<b>Prof. Myoung-Seob Lim</b>	Chonbuk National University, Korea (South)
<b>Prof. Jaime Lloret Mauri</b>	Polytechnic University of Valencia, Spain
<b>Dr. Sotiris Nikolettseas</b>	CTI/University of Patras, Greece
<b>Prof. Bimal Roy</b>	Indian Statistical Institute, India
<b>Prof. Shaharuddin Salleh</b>	University Technology Malaysia, Malaysia
<b>Dr. Lingyang Song</b>	Philips Research, Cambridge, UK
<b>Prof. Guoliang Xing</b>	Michigan State University, USA
<b>Dr. Hassan Yaghoobi</b>	Mobile Wireless Group, Intel Corporation, USA

---

## Editorial Assistants

<b>Shirley Song</b>	Scientific Research Publishing. Email: <a href="mailto:wsn@scirp.org">wsn@scirp.org</a>
<b>Qingchun YU</b>	Scientific Research Publishing. Email: <a href="mailto:wsn@scirp.org">wsn@scirp.org</a>

## TABLE OF CONTENTS

**Volume 1 Number 3**

**October 2009**

**A Cognitive Radio Receiver Supporting Wide-Band Sensing**

V. BLASCHKE, T. RENK, F. K. JONDRAL..... 123

**A Caching Scheme for Session Setup in IMS Network**

Y. F. CAO, J. X. LIAO, Q. QI, X. M. ZHU..... 132

**Metrics and Algorithms for Scheduling of Data Dissemination in Mesh Units Assisted Vehicular Networks**

Z. Y. LIU, B. LIU, W. YAN..... 142

**High Resolution MIMO-HFSWR Radar Using Sparse Frequency Waveforms**

G. H. WANG, Y. L. LU..... 152

**ContSteg: Contourlet-Based Steganography Method**

H. SAJEDI, M. JAMZAD..... 163

**Research on DOA Estimation of Multi-Component LFM Signals Based on the FRFT**

H. T. QU, R. H. WANG, W. QU, P. ZHAO..... 171

**Novel Rate-Control Algorithm Based on TM5 Framework**

Z. J. ZHU, Y. Q. BAI, Z. Y. DUAN, F. LIANG..... 182

**Generation of Multiple Weights in the Opportunistic Beamforming Systems**

G. Y. LU, L. ZHANG, H. Q. YU, C. SHAO..... 189

**The Effect of Notch Filter on RFI Suppression**

W. G. CHANG, J. Y. LI, X. Y. LI..... 196

**Reconfigure ZigBee Network Based on System Design**

Y. XU, S. B. QIU, M. HOU..... 206

**Optimal Deployment with Self-Healing Movement Algorithm for Particular Region in Wireless Sensor Network**

F. ZHU, H. L. LIU, S. G. LIU, J. ZHAN..... 212

**An Adaptive Data Aggregation Algorithm in Wireless Sensor Network with Bursty Source**

K. PADMANABH, S. K. VUPPALA..... 222

# **Wireless Sensor Network (WSN)**

## **Journal Information**

### **SUBSCRIPTIONS**

The *Wireless Sensor Network* (Online at Scientific Research Publishing, [www.SciRP.org](http://www.SciRP.org)) is published monthly by Scientific Research Publishing, Inc., USA.

E-mail: [service@scirp.org](mailto:service@scirp.org)

#### **Subscription rates: Volume 1 2009**

Print: \$50 per copy.

Electronic: free, available on [www.SciRP.org](http://www.SciRP.org).

To subscribe, please contact Journals Subscriptions Department, E-mail: [service@scirp.org](mailto:service@scirp.org)

**Sample copies:** If you are interested in subscribing, you may obtain a free sample copy by contacting Scientific Research Publishing, Inc at the above address.

### **SERVICES**

#### **Advertisements**

Advertisement Sales Department, E-mail: [service@scirp.org](mailto:service@scirp.org)

#### **Reprints (minimum quantity 100 copies)**

Reprints Co-ordinator, Scientific Research Publishing, Inc., USA.

E-mail: [service@scirp.org](mailto:service@scirp.org)

### **COPYRIGHT**

Copyright© 2009 Scientific Research Publishing, Inc.

All Rights Reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as described below, without the permission in writing of the Publisher.

Copying of articles is not permitted except for personal and internal use, to the extent permitted by national copyright law, or under the terms of a license issued by the national Reproduction Rights Organization.

Requests for permission for other kinds of copying, such as copying for general distribution, for advertising or promotional purposes, for creating new collective works or for resale, and other enquiries should be addressed to the Publisher.

Statements and opinions expressed in the articles and communications are those of the individual contributors and not the statements and opinion of Scientific Research Publishing, Inc. We assumes no responsibility or liability for any damage or injury to persons or property arising out of the use of any materials, instructions, methods or ideas contained herein. We expressly disclaim any implied warranties of merchantability or fitness for a particular purpose. If expert assistance is required, the services of a competent professional person should be sought.

### **PRODUCTION INFORMATION**

For manuscripts that have been accepted for publication, please contact:

E-mail: [wsn@scirp.org](mailto:wsn@scirp.org)



# A Cognitive Radio Receiver Supporting Wide-Band Sensing<sup>\*</sup>

Volker BLASCHKE, Tobias RENK, Friedrich K. JONDRAL

*Institut für Nachrichtentechnik, Universität Karlsruhe (TH), Karlsruhe, Germany*

*E-mail: {blaschke, renk, fj}@int.uni-karlsruhe.de*

*Received April 29, 2009; revised May 8, 2009; accepted May 10, 2009*

## Abstract

The specification of IEEE 802.22 defines the world-wide first cognitive radio (CR) standard. Within a range of 40 MHz to 910 MHz CR systems are allowed to allocate spectrum besides the currently established radio services like radio and TV broadcasting. In order to fulfill the regulative guidelines of interference limitations, a capable spectral sensing and user detection has to be provided. Due to the wide frequency range specified in IEEE 802.22 and the high dynamic range of signals allocated in this band there are high demands on the CR receiver's front-end. Especially the performance requirements on analog-to-digital converters increase significantly compared to current wireless systems. Based on measurements taken in this frequency range requirements to CR's ADCs are figured out. Furthermore, the measurement results are analyzed regarding expectable allocation scenarios and their impacts to spectral sensing. Derived from this results and a comparison of general spectral sensing mechanisms an approach for a CR receiver supporting wide-band sensing is presented. Considering the apriori information resulting from scenario analysis and including adapted information processing in the terminal the ADC's performance requirements can be reduced.

**Keywords:** Cognitive Radio, IEEE 802.22, Spectrum Sensing, A/D Conversion

## 1. Introduction

The term *Spectrum Sensing* becomes more and more important, especially in the context of cognitive radio (CR). Due to the increased request for wireless transmission resources and the ongoing installation of new radio access technologies for broadband access, enhanced research in the field of mobile CR receivers is necessary. Based on the results of spectral measurements [1,2] a low utilization over wide frequency ranges was identified. This additionally motivates the development of intelligent radio resource allocation mechanisms to overcome this waste of resources. For increasing the overall utilization, dynamic allocation of free spectral resources that considers both the users and the spectral environment is required. This approach is supported by the CR concept [3,4]. Providing a dynamic resource allocation, sufficient information about the spectral environment has to be collected. Thereby, different acquisition methods

can be used. On one side, all information are collected by a central control unit and distributed to simple mobile terminals. In this case traffic load information could be exchanged between joint networks via backbone [5]. Other approaches base on distributed sensing using all mobile terminals of a radio access network. This requires sensing capabilities in each mobile entity and an efficient algorithm for consolidation of the results. Especially in such scenarios swarm-intelligence algorithms could offer additional benefits. Nevertheless, an appropriate spectral sensing and information extraction forms a precondition for dynamic and efficient allocation mechanisms. In order to avoid unacceptable interferences a reliable detection of other users has to be supported by the CR. Therefore, both, the temporal as well as the spectral characteristics of the observed frequency band have to be known to the terminal. All these requirements lead to high demands on the radio's frontend.

In this paper the performance demands on analog-to-digital converters (ADCs) in mobile CR receivers supporting spectral sensing are discussed. Based on the frequency ranges specified in IEEE 802.22 [6] general as-

<sup>\*</sup>This paper is an expanded version of the correspondent article accepted in the proceedings of "2008 IEEE International Conference on Communications Workshop."

pects of wide-band sensing as well as specific demands on the structure of mobile CRs are presented. Analyzing the frequency range of IEEE 802.22 several types of channel utilization can be figured out. This a-priori knowledge combined with suitable information processing in the terminal will lead to an optimized front-end structure supporting wide-band spectral sensing.

The paper is structured as follows: In the next section a brief introduction to ADCs is presented. In Section 3, a detailed description of sensing algorithms in CR terminals is given. Considering the frequency bands specified in IEEE 802.22 the expected signal characteristics are identified. Based on this, two general wide-band sensing methods including a comparison of their demands to ADC's performance are described in Section 4. In Section 5 a CR receiver structure is presented including adapted spectral sensing combined with a convenient information processing. Finally, a conclusion is given.

## 2. Analog-to-Digital Conversion in CR Receiver

The digitization of the received signal is a basic component in each digital receiver. Only a suitable sampling and quantization of the analog input signal enables the receiver to provide the communication tasks supported by digital signal processing. Most CR concepts described in literature assume an appropriate ADC as precondition. But an efficient analog-to-digital conversion contains a lot of challenges in order to support the performance constraints assumed in these CR concepts.

Generally, the key parameters for summarizing the ADC's performance are stated resolution, signal to noise ratio, spurious free dynamic range, and power dissipation [7]. Furthermore, aperture jitter as well as two-tone intermodulation distortion is important for characterizing ADCs. A detailed description and performance analyzes can be found in [8] and [9]. In [8] the performance of on-the-market ADCs is analyzed in order to describe the evolution and trends in ADC's technology. This evaluation was continued in [9] including also present-day trends.

In this article ADC parameters which restrict an implementation of current ADCs in mobile CR terminals are figured out. These are sampling frequency  $f_s$ , affecting the effective resolution bandwidth, effective number of bits  $N_{\text{eff}}$ , describing the dynamic range supported by the ADC, and power dissipation  $P_{\text{diss}}$  resulting in battery running time.

In order to fulfill the Nyquist criterion the converter's sampling frequency  $f_s$  has to be more than two times the effective analog bandwidth [10]. Furthermore, the effective number of bits  $N_{\text{eff}}$  is lower than the stated number of resolution bits specified by the vendor. Due to hardware imperfections and quantization noise  $N_{\text{eff}}$  is [8]

$$N_{\text{eff}} = \frac{D - 1.76}{6.02} \quad (1)$$

where  $D$  describes the effective dynamic range of the converter in dB. Especially for detection and sensing applications a high dynamic range is of increased importance. If the dynamic range of the expected input signals is higher than ADC's dynamic range, weak signals may not be detected due to resolution limitations.

Having a look to the results depicted in [8] and [9] three main hardware architecture concepts become potential candidates for implementation on mobile CR terminals. Flash converters offer sampling rates of about 1 Gsps due to a parallel comparator structure. But this requires high hardware effort which causes increased power dissemination. Therefore, this architecture is unattractive to mobile applications. Due to the high hardware effort an effective resolution of only  $N_{\text{eff}} = 6 \dots 8$  Bits can be supported. Using Pipelined ADCs  $N_{\text{eff}}$  can be increased up to 15 Bits but the sampling frequency  $f_s$  is less than 500 Msps. Due to the specific design, implementation of analog track-and-hold blocks is required. The third group of potential candidates is  $\Sigma\Delta$ -converters. The available effective resolution is up to 20 Bit but the available sampling rate is less than 100 Msps. Though, due to their low power consumption these ADCs are very interesting for an implementation in mobile CR terminals. A detailed description of the different ADC structures can be found in [10,11].

## 3. Spectrum Sensing and Related Hardware Impacts

Sensing the current channel state is one important task of each radio system using dynamic channel allocation. The well-known IEEE 802.11 wireless LAN systems, for instance, use carrier sense multiple access with collision avoidance (CSMA/CA). This multiple channel access scheme requires sensing and detection of channel allocation. In such systems, sensing bandwidth and transmission channel bandwidth are the same.

During the last years several approaches for dynamic spectrum allocation (DSA) [12] and channel allocation adapted to the current spectrum situation and user requirements were published [5]. These concepts combine available transmission resources of several systems in order to optimize the spectrum utilization of all considered systems. This would require the availability of the current system state to all other systems. Assuming a coupling of combined systems and a general control entity, the information exchange can be realized using traffic control channels in the wired backbone network. So, each system can handle and optimize the channel allocation of its subscribers considering this additional traffic

load information. Inter-system handover need to be initiated and controlled by the general control entity. Basically, this architecture requires cooperation between all radio access networks participating in DSA, which again results in a relatively static system configuration.

Another approach for increasing the spectral utilization is overlay systems [13]. Based on the fact that wide frequency ranges offer a lot of unallocated transmission capacity, this approach describes the allocation of local networks exploiting temporarily and/or locally unoccupied transmission channels. Due to the basic precondition that licensed systems should not be changed for overlay usage, the rental user must observe the communication channel in order to provide a reliable detection of the licensed user's channel allocation.

Furthermore, the rental user's signal must be adapted to the licensed system's transmission parameters regarding channel bandwidth, maximum channel allocation duration, transmission power, etc. Static system parameters, e.g., channel bandwidth, can be defined in a database available for each rental user. But for dynamic parameters, e.g., current channel allocation, licensed user's allocation statistics have to be observed and analyzed at present. So, a continuous spectral observation has to be done by the overlay system. In order to get information about the allocation of frequency bands, energy detection can be used. Comparing the received signal power with the noise level general channel allocation information can be collected. If the signal power is higher than the measured noise level the channel is already occupied. Especially, in case of weak signals that are close to noise level the power detection can fail. So, analyzing signal's higher order statistics or other feature detectors may overcome this drawback [14].

Having a look to the IEEE 802.22 specification, the communication channel bandwidth is  $B_{ch} = 6 \dots 8$  MHz and the frequency range specified for allocation is between 41 MHz and 910 MHz with respect to national regulations [6]. So, the overall system bandwidth, that has to be observed, is  $B_s = 869$  MHz. This is more the 100 times the signal's bandwidth  $B_{ch}$ . In order to provide a flexible CR system which is able to optimize spectral utilization, the full frequency range has to be supported by mobile terminals. This also includes the ability for a fast and efficient sensing of wide ranges in order to adapt transmission parameters to the licensed user's allocation. Supporting high signal bandwidth directly affects the hardware architecture of a terminal. Especially the analog signal processing and the ADC limit the supported bandwidth. As it is described in Section 2, there is a trade-off between the signal's bandwidth and the dynamic range of the converter.

Having a look to the receiving signals within this frequency range, several allocation characteristics can be

pointed out. In Figure 1 the received signal power per frequency averaged over a sensing period of 3 h is depicted. The frequency range of  $f = 41 \dots 910$  MHz represents the overall range specified for IEEE 802.22. As it can easily be seen, there are wide frequency ranges where a very low averaged signal energy is detected. But also high utilized bands can be pointed out. Between 88 MHz and 108 MHz the European FM radio broadcast service is allocated (cf. ch 1). Furthermore, some TV broadcast signals as well as temporarily allocated channels can be noticed. For describing the channel utilization during sensing time a binary spectrogram can be defined:

$$O\{x(t)\}_{|f=f_m} = \begin{cases} 1 & \text{for } S\{x(t)\}_{|f=f_m} > P_{th} \\ 0 & \text{for } S\{x(t)\}_{|f=f_m} \leq P_{th} \end{cases}, \quad (2)$$

where  $x(t)$  is the received signal,  $S\{x(t)\}_{|f=f_m}$  describes the spectrogram of  $x(t)$  at frequency  $f_m$ , and  $P_{th}$  describes the detection threshold. The resulting binary description of the spectral allocation can be averaged over sensing time using a window length  $N_w$ . So, the averaged channel utilization can be written as

$$\bar{O}_{t_n}|_{f=f_m} = \frac{1}{N_w} \sum_{t=t_n}^{t_n+N_w+1} O_{t_n}|_{f=f_m}, \quad (3)$$

where  $t_n$  is the time index. The resulting characteristics for three different frequencies are depicted in Figure 2. The curve of channel 1 describes a typical broadcast channel utilization. Channel 2 offers a varying averaged occupation between 0.01 and 0.92. In channel 3 an averaged utilization of 0 can be noticed. During this period the channel is not used by the licensed user and would be interesting for CR resource allocation. In order to detect

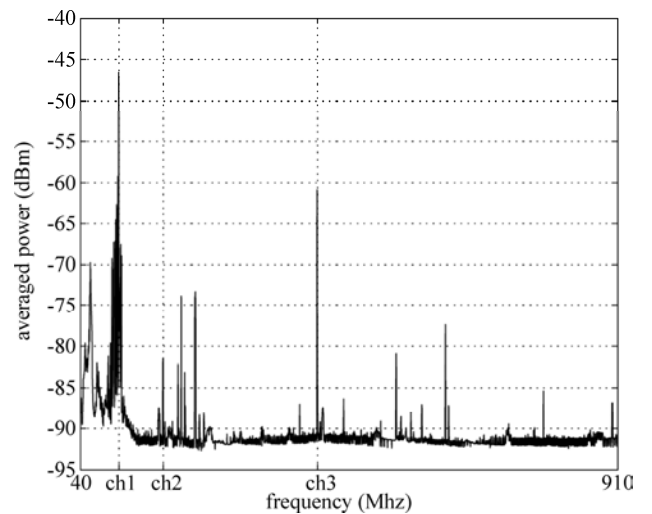


Figure 1. Averaged power vs. Frequency in IEEE 802.22 frequency range.

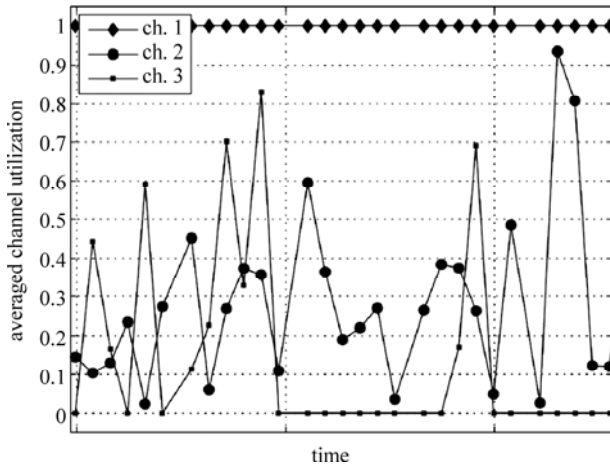


Figure 2. Averaged channel utilization of channel 1, 2 and 3.

such periods that are also called white spaces [4] a suitable observation of wide frequency ranges has to be provided. As it can be noticed in Figure 1, high differences in the dynamic range of the depicted signals has to be considered. During measurements a maximum signal power of  $-41.8$  dBm at channel 1 was observed. The general noise level is  $-92$  dBm measured in unallocated sub-bands. So, the overall dynamic range compared to the bandwidth of more than 850 MHz marks the main challenge in finding suitable solutions for mobile CR receivers supporting this wide frequency range. In the next section two general sensing methods supporting these requirements will be described.

#### 4. Spectrum Sensing Methods

Several methods can be used for analyzing wide ranges of radio spectrum. The two general concepts, sweeping a small detection window over the observed frequency band and wide band analog to digital conversion followed by energy or feature detection are briefly de-

scribed in the following subsections. Their ability for implementation in a mobile CR receiver will also be discussed.

##### 4.1. Basics on Spectral Analysis

Due to the fact that currently available ADCs possess a limited bandwidth, it is not possible to digitize large frequency spans at once. Therefore, so-called fast Fourier transform (FFT) analyzers are only usable for low frequency signals. In order to investigate high frequency signals, superheterodyne receivers must be applied. Here, the overall input frequency span is mixed to a common intermediate frequency (IF) by a tunable local oscillator [15]. Spectral resolution is directly determined by the IF filter. The smaller the resolution bandwidth  $B_R$ , the higher the spectral resolution. A well-known problem occurs if the input frequency range is more than two times bigger than the IF, because then suppression of the image frequency is not possible without affecting the input signal. Hence, a tunable bandpass is necessary. This problem can be overcome if several IF stages are used, where the first one transforms the input signal to a higher frequency. It is then possible to suppress the image frequency without affecting the input signal.

Figure 3 shows a superheterodyne receiver with two IF chains for the IEEE 802.22 specification. The RF input signal first passes RF attenuation that helps to prevent overload and distortion. Afterwards, a preselector lowpass filters out higher frequency signals. In order to mix the RF signal up to 1000 MHz ( $f_{IF1}$ ), the first local oscillator (LO 1) must operate in a frequency range from 1041 MHz to 1910 MHz. This leads to an image frequency that ranges from 2041 MHz to 2910 MHz which can easily be suppressed by the following IF 1 filter. After this, a second local oscillator (LO 2) with a frequency of 970 MHz mixes the 1000 MHz signal down to  $f_{IF2} = 30$  MHz. Image frequency is 940 MHz which

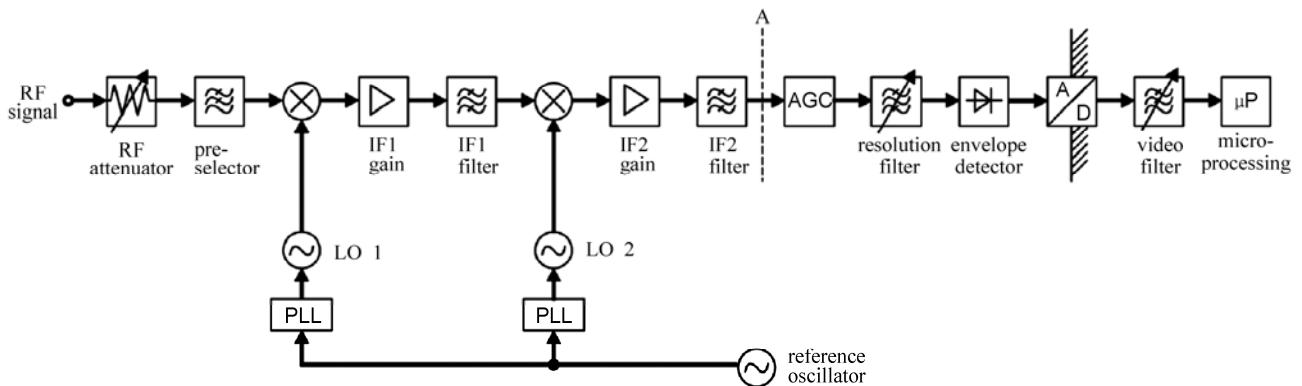


Figure 3. Superheterodyne receiver with two IF chains and low frequency ADC for the IEEE 802.22 specification.

again can easily be filtered out by the IF 2 filter. Both LOs are controlled by PLLs that are connected to a reference oscillator to increase frequency accuracy. The automatic gain control (AGC) block is followed by a tunable bandpass filter that determines the resolution bandwidth. The ideal case for the resolution filter is a rectangular filter with bandwidth  $B_R$ . In order to achieve short measurement times, however, optimized Gaussian filters are used that are temperature stable and possess a higher bandwidth accuracy as well. Nevertheless,  $B_R$  influences the sweep time  $T_{sw}$  that is necessary to scan the whole frequency range  $B_S$ . If measurement time falls below  $T_{sw}$ , amplitude losses and signal distortions occur that eventually lead to frequency offsets. Subsequent to filtering is the envelope detector and ADC. The video filter is a lowpass that suppresses noise and helps to smooth the signal spectrum. The video bandwidth  $B_V$  acts inversely proportional to the sweep time  $T_{sw}$ . The micro-processing block ( $\mu P$ ) includes, e.g., averaging and threshold decision making.

For the definition of the required sweep time, two cases must be taken into consideration, one where the video bandwidth is higher than the resolution bandwidth and vice versa [15]:

$$T_{sw} = \begin{cases} k \cdot \frac{B_S}{B_R^2} & \text{for } B_V > B_R \\ k \cdot \frac{B_S}{B_R \cdot B_V} & \text{for } B_V < B_R \end{cases}, \quad (4)$$

where the parameter  $k$  denotes a proportional factor that is usually in the range of 1 to 3. The first case of Equation (4) is illustrated in Figure 4. System parameters were chosen according to the IEEE 802.22 specification. The vertical dashed lines define the channel bandwidth of  $B_{ch} = 6 \dots 8$  MHz [16]. It can easily be seen that for one

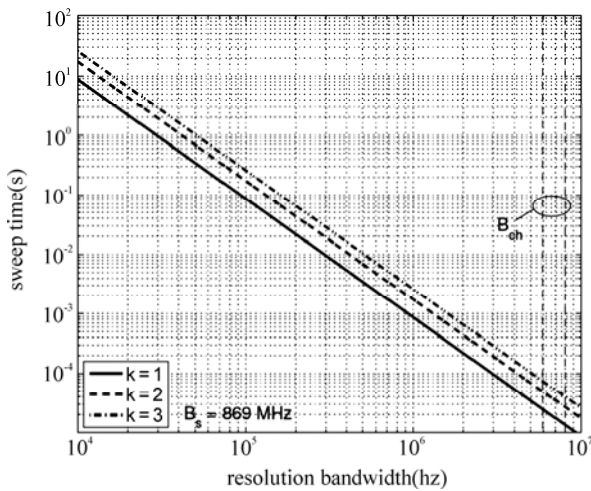


Figure 4. Sweep time  $T_{sw}$  in relation to resolution bandwidth  $B_R$  with parameter  $k=1,2,3$ .

single scan through the whole frequency range of 869 MHz approximately  $3 \cdot 10^{-4}$  s are necessary if we consider a resolution bandwidth of  $B_R = 2.5$  MHz. (This means that we take two to four frequency bins per channel.)

As it can be seen in Figure 3, the ADC is located after the resolution filter and the envelope detector. This means that only a bandwidth of  $B_R$  has to be digitized which leads to a higher amplitude resolution. Moreover, a quite simple ADC is sufficient in that case. Drawbacks are, however, that pretty much analog hardware is necessary and that the measurement time is increased (cf. Section 5). Additionally, no feature detection is possible since only energy detection is performed.

## 4.2. Wide-Band Sensing

Another possibility for energy detection is the direct conversion of the wide-band input signal. This reduces the number of intermediate frequency stages required for sweeping the detection window but significantly increases the performance requirements at the ADC.

Ideally, the incoming analog signal is bandpass filtered by a bandpass with the system bandwidth  $B_S$ . Afterwards, the signal is amplified and down converted from the radio frequency  $f_{RF}$  to an intermediate frequency  $f_{IF}$ . Subsequent to a second filtering and AGC, the analog to digital conversion is done. The following data processing includes an FFT in order to extract the current power allocation over the frequency.

After digitization the signal contains information of the complete observed frequency range. Of course, the information depth is characterized by the resolution performance of the ADC, which is a main drawback of this approach. As it was shown in the section before, the observed signal is characterized by a high variation of the spectral power density. The overall dynamic range is more than 50 dB. Furthermore, under-utilized small-band signals with weak signal amplitudes can be noticed. In order to provide a reliable detection of the licensed user, these signals still have to be noticeable after the ADC, otherwise the signal detection fails. Therefore, the dynamic range of the ADC is an important parameter for CR terminals.

Based on Equation (1) the presented measurements would require a minimum resolution of  $N_{eff} = 9$  Bits. Besides a high bit resolution, also a high sampling rate is required for a wide band digitization. Following the example of IEEE 802.22, the sampling rate is  $f_{samp} = 1.82$  Gsps meeting the Nyquist criterion. Due to the overall frequency range from 41 MHz up to 910 MHz, a bandpass sub-sampling cannot be used for reduction of  $f_{samp}$ .

As stated in [9], the maximum sampling rate for an effective resolution of  $N_{eff} = 9$  Bits is about  $f_{samp} = 500$  Msps ... 1 Gsps. The analysis in [8] and [9] also figure

out those sampling rates significantly higher than 100 Msps can only be handled by Flash or Pipelined converters. As described in Section 2, these two architectures are characterized by a high power consumption and, therefore, not preferable for an implementation in mobile terminals. The group of  $\Sigma\Delta$ -converters offer lower power consumption but cannot provide the high sampling rate available for Flash converter. In order to suppress the undesired signals also analog notch filter can be applied. This requires a first scan for identification of the strongest signals. After tuning the notch filter to these frequencies a second scan is used for detection of weak signals. Besides doubling the scan time also the analog hardware effort increases significantly.

Compared to the energy detection described above, a wide band digitization offers a better time resolution in a wide frequency band. Thus, a detailed extraction of temporal features is possible. A high time resolution becomes important for detection of very short channel allocations or for detailed analyzing of the licensed user's allocation statistics.

Both concepts offer different advantages which are required for a successful operation of flexible overlay systems. On the other side, each concept has significant drawbacks that can not be solved within the next couple of years and preclude an implementation in mobile terminals. So, a combination of both could give the opportunity to combine the advantages. This approach is discussed in Section 5.

## 5. Wide-Band CR Receiver for IEEE 802.22

### 5.1. Sub-Band Spectrum Sensing

As we presented above, a digitization of the complete system bandwidth  $B_S$  is not useful regarding to technical and economical constraints. On the other side, the extension of spectral sensing described for CRs will support an increased spectral utilization. In order to simplify the spectrum sensing and to reduce the hardware requirements described in Subsection 4.1., we will analyze the measurement results with respect to signal characteristics expectable in the IEEE 802.22 frequency range. Having a closer look to the results depicted in Figure 1 and Figure 2, it can be seen that channel utilization in the sub-band allocated by broadcast services (cf. ch 1) does not change during the measurement time. Furthermore, the received signal power of broadcast transmitters averaged over the sensing time is significantly higher than the noise level. In general, the average power level of the broadcast signals is higher than  $-70$  dBm and even higher than  $-50$  dBm considering the strongest signal. Therefore, a continuous sensing of these frequency

ranges does not offer any additional information for an operating CR system. This leads to the opportunity to exclude such quasi-static frequency ranges considering its information entropy during the spectrum sensing phase. Furthermore, a decreased dynamic range of the input signal reduces the ADC's hardware recommendations. Without loss of generality statistical independence of observed communication channels can be assumed. So, the spectrum that is observed for detection of averaged channel allocation can be divided into several sub-bands. These sub-bands do not need to be observed at the same time but can be sensed sequentially as long as the sensing is repeated periodically and the sensing interval as well as the sensing period is suitable to the licensed user's signal. Based on this knowledge and assuming a suitable sensing sequence the full system band can be split into  $M$  sub-bands:

$$B_S = M \cdot B_{\text{sub}} \quad (5)$$

Each sub-band is separately digitized reducing the sampling rate of the ADC. Due to decoupling highly utilized communication channels and frequency bands with low spectral utilization, the dynamic range of the ADC's input signal can be optimized as well, which results in an enhanced sensing of weak signals. Additionally, the suppression of strong signals using analog notch filter is not necessary, because the reduced sampling rate enables higher bit resolution.

### 5.2. Receiver Structure

The system architecture described in Subsection 4.1 defines the bases for our wide-band CR supporting sub-band sensing. In contrast to the structure depicted in Figure 3 the ADC is placed directly after the IF 2 filter. This position is marked with the letter 'A'. The adapted signal processing of the wide-band CR receiver is depicted in Figure 5. Until the marker 'A' the analog signal processing is the same as described in Subsection 4.1. In order to support the IEEE 802.22 specifications the sub-band bandwidth is defined to  $B_{\text{sub}} = 50$  MHz including six sub-channels at a bandwidth of 8 MHz up to eight sub-channels at a bandwidth of 6 MHz, respectively. Similar to the number of sweep points defined in analyzer detection the LO 1 can be tuned to 20 predefined frequency steps resulting in a small sub-band overlap. For a more flexible sub-band configuration, a continuous oscillator tuning could also be implemented. Generally, digitization of a 50 MHz sub-band requires a sampling rate of about 100 Msps. In the following digital sensing processing energy detection as well as feature detection or other signal detection algorithms could be applied. In Figure 5 the block structure for energy detection is depicted. Due to the digital processing the single scan time

can be reduced by factor  $1/20$  compared to the analog processing using  $B_R = 2.5$  MHz.

Besides the reduction of scan time the proposed sub-band digitization offers the possibility of sensing adjacent channels of the currently allocated communication channel. As it is depicted in Figure 5 the received signal is used for communications and sensing processing in parallel. This means that all spectral information within the actual sub-band can be collected along with current data transmission. Furthermore, the additional time released by the parallel sensing and communication processing can be used for additional sensing of other sub-bands. The additionally obtained information increases the CR's knowledge about its spectral environment. But besides an efficient sensing algorithm, also a suitable information processing and knowledge storage has to be applied in mobile CR terminals. In the next subsection an algorithm is presented what bases on the proposed sub-band sensing.

### 5.3. Information Processing

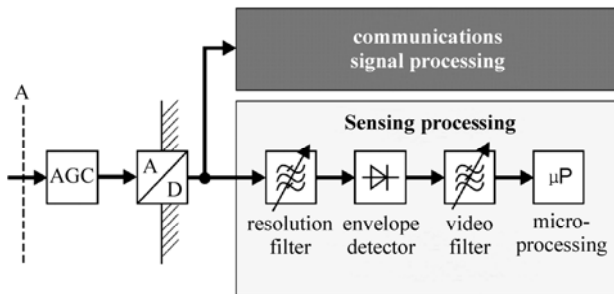
Another question in the context of CRs deals with processing the data gained from spectrum sensing. As it was described by Mitola [3], one important enhancement of CRs compared to SDRs is the implementation of reasoning algorithms. Reasoning can be applied for user centric applications like user interface adaptation or providing user-specific local information and news services. Moreover, processing of experiences can also be used for optimizing the spectral sensing procedure. As long as a predefined performance can be fulfilled, it is not necessary to change the terminal's configuration or the allocated communication channel. Hence, sensing the full spectral range is not necessary for most of the time. In case that the currently occupied sub-band cannot longer be allocated by rental users the information of the next most suitable sub-band is required. Due to the possibility of a direct sub-band scanning described above the scanning procedure can be limited to such sub-bands which offered low channel utilization in the past. In order to get

a first overview of a sub-band, the averaged utilization (cf. (3)) of a communication channel or a complete sub-band could be considered. Due to the digitization of a full sub-band  $B_{\text{sub}}$ , all communication channels in this band can be observed simultaneously. As long as the current sub-band can offer some free radio resources to the CR overlay system, other sub-bands need not be observed continuously. A periodical short scan provides information to approximately trace the sub-band utilization. Based on this sensing result, the sub-band can be ordered with respect to the current utilization. In case of shifting overlay users to another sub-band, the sub-band offering the lowest utilization can be observed in detail.

Since spectrum sensing is only one task of a CR terminal, besides radio communication, an intelligent scheduling of the sensing periods is necessary. In case of a feasible number of CR nodes at one location, the quality of the detection result will not increase significantly compared to the number of additional nodes. Thus, distributed sensing of different sub-bands that is provided by several nodes at the same time will help gain more information of the overall frequency range. In order to use this advantage, the challenge of collecting the information from all distributed nodes need to be solved.

In [13] an innovative approach for distributed sensing that provides a solution to overcome the hidden-node problem was proposed. In the described system all nodes sense the same sub-band synchronously. After the sensing, all binary detection results are collected at one master station of the local network. The detection results of each single channel within the observed sub-band are coded in a one bit decision. These bits are sent simultaneously by all nodes that are connected to the master station. For a detailed description of the algorithm the reader may be referred to [13]. Basically, the simultaneous transmission leads to a superposition of all detection results, which can be interpreted as a logical OR operation. Thus, a reliable detection of an increased area can be provided.

Adopting this approach to the problem of increasing the observed spectral range, we can use the following strategy: Besides the signaling of the detection results within the currently used sub-band another sensing and signaling period can be added. During this period an additional sensing and signaling of adjacent sub-bands may be executed. Following the signaling method described in [13] the sensing results are superposed. Thus, the calculated sub-band utilization is a rough estimation of the available resources. In case of a low utilization a detailed sensing will follow. The measurement process supporting such a distributed sensing is depicted in Figure 6. The data transmission including sensing and signaling regarding [13] in the currently used sub-band is named *S1*. During this phase dedicated nodes observe also adjacent sub-bands. These sensing results are sent to the



**Figure 5. Structure of wide-band CR receiver supporting sub-band sensing.**





**Figure 6. Sensing process for wide-band CRs supporting sub-band processing.**

central control station during the phase *S2*. Furthermore, a full range scan could be initiated, which is named *S3*. Since the full scan needs, however, more time compared to the normal periodic sub-band scan, the communication in the CR system may be affected by the full scan. Due to the ranking of the sub-band utilization, only the most suitable sub-bands need to be considered. A full processing of allocation information is only required for the current sub-band. All other sub-bands are characterized by an averaged utilization index that reduces the memory and processing effort in the mobile terminal. For a detailed investigation of the proposed distributed sensing method, e.g., bio-inspired algorithms could be taken into consideration.

## 6. Conclusions

In order to increase spectral utilization future CR receivers have to provide spectrum sensing capabilities. Applying DSA mechanisms require suitable spectrum sensing capabilities in order to adapt the radio transmission to the identified spectral environment. In the CR standard IEEE 802.22 a frequency range from 41 MHz to 910 MHz is specified. Within this spectral band CR networks can be established under the limitation that the already established radio services are not interfered unacceptably. In order to provide a reliable signal sensing and detection in the CR terminals several preconditions to the receiver's front-end have to be fulfilled. In this paper demands on ADCs in such wide-band scenarios are discussed in detail. As presented in Section 2, the general performance of ADCs is characterized by the trade-off between supported bandwidth and dynamic range defined by the effective number of bits. Today's ADC support sampling rates up to several Gsps at the expense of low dynamic range and high power consumption. But the demand for a high sensing quality in mobile receivers leads to the contrary request for high dynamic ranges at a low power consumption. Thus, the input bandwidth has to be reduced. In Section 3 the different radio services allocated in the considered frequency range are analyzed. It is shown, that the utilization varies significantly over the frequency. Thus, different demand for sensing in the sub-bands can be observed. Especially sub-bands allocated by broadcast radio services do not offer additional radio transmission resources but increase the demands on ADC's performance significantly. Based

on the sensing algorithms described in Section 4, a receiver for wide-band mobile CRs based on the superheterodyne principle is presented in Section 5. Due to the proposed distribution of the wide-band into 20 sub-bands high-utilized sub-bands can be masked while sensing of under-utilized sub-bands benefits from the increased resolution of the ADC resulting from the decreased input signal bandwidth. It combines this sub-band sensing method with the presented information processing results in a capable mobile receiver structure for IEEE 802.22 CR networks.

## 7. References

- [1] FCC, "Spectrum policy task force report, ET Docket No. 02-155," Technical Report Series, November 2002.
- [2] Shared Spectrum Company, "Comprehensive spectrum occupancy measurements over six different locations," August 2005, <http://www.sharespectrum.com/>.
- [3] J. Mitola, "Cognitive radio-an integrated agent architecture for software defined radio," Ph.D. dissertation, Royal Institute of Technology (KTH), Kista, Sweden, 2000.
- [4] S. Haykin, "Cognitive radio: Brain-empowered wireless communications," *Selected Areas in Communications*, IEEE Journal on, Vol. 23, No. 2, pp. 201–220, February 2005.
- [5] P. Cordier *et al.*, "E2R cognitive pilot channel concept," IST Mobile Summit, Mykonos, Greece, June 2006.
- [6] IEEE, "P802.22: Cognitive radio, wide regional area network," Technical Specifications, May 2005.
- [7] R. H. Walden, "Analog-to-digital converter technology comparison," *IEEE GaAs IC Symposium Technical Digest*, pp. 217–219, October 1994.
- [8] R. H. Walden, "Analog-to-digital converter survey and analysis," *Selected Areas in Communications*, IEEE Journal on, Vol. 17, No. 4, pp. 539–550, April 1999.
- [9] B. Le, T. W. Rondeau, J. H. Reed, and C. W. Bostian, "Analog-to-digital converters," *Signal Processing Magazine*, IEEE, Vol. 22, No. 6, pp. 69–77, November 2005.
- [10] R. Plassche, *CMOS Integrated Analog-to-digital and Digital-to-analog Converters*, 2nd Edition, Kluwer Academic Publishers, Boston/ Dordrecht/London, 2003.
- [11] B. Brannon, *Software Defined Radio-Enabling Technology*, John Wiley and Sons, London, W. Tuttlebee, Ed., ch.

- Data Conversion in Software Defined Radios, pp. 99–126. 2002.
- [12] P. Leaves *et al.*, “A summary of dynamic spectrum allocation results from drive,” in IST Mobile and Wireless Telecommunications Summit, pp. 245–250, June 2002.
- [13] T. A. Weiss and F. Jondral, “Spectrum pooling: An innovative strategy for the enhancement of spectrum efficiency,” *Communications Magazine*, IEEE, Vol. 42, No. 3, pp. 8–14, March 2004.
- [14] D. Cabric *et al.*, “Implementation issues in spectrum sensing for cognitive radios,” in *Signals, Systems and Computers*, 2004, Conference Record of the Thirty-Eighth Asilomar Conference on, Vol. 1, pp. 772–776, November 2004.
- [15] C. Rauscher, *Grundlagen der Spektrumanalyse*, Rohde & Schwarz, 2004.
- [16] C. Cordeiro *et al.*, “IEEE 802.22: An introduction to the first wireless standard based on cognitive radios,” *Journal of Communications*, Vol. 1, No. 1, April 2006.

# A Caching Scheme for Session Setup in IMS Network

Yufei CAO<sup>1,2</sup>, Jianxin LIAO<sup>1,2</sup>, Qi QI<sup>1,2</sup>, Xiaomin ZHU<sup>1,2</sup>

<sup>1</sup>State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications,  
Beijing, China

<sup>2</sup>EBUPT Information Technology Co., Ltd, Beijing, China

E-mail: {caoyufei, liaojianxin, qiqi, zhuxiaomin}@ebupt.com

Received April 29, 2009; revised May 8, 2009; accepted May 10, 2009

## Abstract

In IP Multimedia Subsystem (IMS), the session setup delay is a critical value for Quality of Service (QoS). The existing approaches to improve this metric depend on optimization of Session Initiation Protocol (SIP) message transmitting and signaling flows. Unfortunately, some service features are missing considered although they have been used widely in traditional 2G networks. This paper proposes a novel session setup scheme based on caching, upon the investigation of the performance of IMS session establishment. This mechanism involves cache based local routing policy and an adaptive caching algorithm, which can decrease call setup delay effectively as cached information in the terminating Serving-Call Session Control Function (S-CSCF) hit. The analytical model is deduced, as well as the delay and cost ratio functions are presented based on the model. Moreover, the analytical model is validated through the performance simulation in which the performance of the proposed novel method is evaluated against the basic session setup mechanism in terms of cost and delay.

**Keywords:** SIP Session, Caching, IMS

## 1. Introduction

IMS (IP Multimedia Subsystem) network is introduced in 3GPP R5, which aims to provide mobile user multimedia services such as voice, video and data. It unifies core network as all-IP network architecture and realizes the integration of fixed and mobile communication networks [1–3]. Home service control is selected in IMS, which means the entity that accesses to the subscriber database and interacts directly with service platforms is always located at the user's home network. Thus, location management and session management are pointed to the home network as far as possible. The HSS (Home Subscriber Server) contains all the information related to the users and their services. The S-CSCF (Serving-Call Session Control Function) located in home network provides session control and registration services [2,4,5].

In IMS network, when the caller A wants to establish session with the callee B, SIP INVITE request constructed by UE (User Equipment) is forwarded to the user A's home network via the P-CSCF (Proxy-Call Session Control Function). And then user A's home S-CSCF executes the service control, including interaction with the AS (ap-

plication server), a process of querying DNS to determine the entry of UE B's home network and assigning the S-CSCF through the I-CSCF (Interrogating-Call Session Control Function) which is needed to select the S-CSCF of UE B. This S-CSCF is responsible for dealing with and ending the session, containing the interaction with AS, sending request messages to the P-CSCF which UE B accessed to and forwarding to UE B finally. The response generated by UE B reverses the same path back to UE A. After several forward and back flows, session establishment is completed [6,7].

The excessive signaling of current IMS session setup mechanism results in the long delay from the session initiation of caller and the final response of callee [8,9]. This is unfavorable to those applications which require fast communication handshake. Comparing with traditional 2G network, there are some problems in the flows defined in current specifications: 1) There is no considering that users' session setup takes place within one S-CSCF serving area, so DNS querying and S-CSCF assignment are involved. Thus, this process brings two problems. Firstly, it increases unnecessary interactive signaling and the session setup delay is raised. Secondly, the larger the scale of

users is, the more the traffic load of the entities as DNS, I-CSCF, HSS, etc. is. That not only would be a waste of network bandwidth and resources, but also reduces system reliability. And eventually it leads to long session setup delay for users' experience. 2) A series of interactive inquiries to establish sessions is too complicated. Even if the requested destination server addresses (S-CSCF address) of the current session request is same as the previous one, DNS inquiry process and the S-CSCF selection are still necessary. So it is feasible to optimize signaling traffic load by decreasing the number of signaling interaction.

This paper focuses on the study of basic session setup mechanism in current IMS core network, which includes the process of caller's session request arriving at its S-CSCF (the originating S-CSCF) and the connection setup between the originating S-CSCF and the terminating S-CSCF. A cache based session setup mechanism is proposed by improving the originating S-CSCF session setup procedure with taking locality and caching into account. The advantages are: 1) when session participants are in the one S-CSCF, their sessions can be established directly. 2) If the destination server address of current session request is the same as the previous one, session can be established directly. It decreases the times of signaling interaction for DNS query and reduces the load of the HSS by improving the S-CSCF assignment process. 3) There is no change of IMS core network, no adding or varying to the terminal signaling, and no impact on signaling flows. The simulation shows that the cache based method is able to reduce signaling traffic load and session setup delay in IMS. At the same time, as the improvement is mainly about the signaling flows, instead of system hardware or network structure, the cost of its implements is smaller comparatively.

The rest of the paper is organized as follows. Section 2 reviews related work and motivates the cache based session setup in IMS. Section 3 presents the basic session setup flows and analyzes the terminating S-CSCF routing process with explanation of its issues. Section 4 introduces the cache based session setup mechanism with original contribution presented in detail. Section 5 shows the detailed analysis of the cost and the mean delay function of new mechanism. Section 6 shows evaluation of the performance by simulation. Section 7 concludes this paper.

## 2. Related Work

Usually, in Circuit Switch (CS) of both traditional 2G and 3G networks, there are two approaches to improve call setup performance [10].

One is to improve location management policy and management protocols [11–14], which aims at exploring how to manage and query user's location information

efficiently, in order to quickly address entries serving the callee in the core network. For example, there are three-tier location management in [11], and layered management of mobile IP in [12,13] which restricts UE register signaling in its local networks. Also, in [14], a cache scheme of location information is proposed to reduce call setup delay. In [3], two functional entities MMS (mobile management server) and ACS (access control server) are introduced to enhance 3G core network architecture. It separates registration procedure and security service away from CSCFs to achieve the efficiency of session setup by allaying complexity of CSCFs without change to the current IMS session setup procedure.

The other one is to improve call setup process by taking advantage of the users' locality to advance call setup performance [10,15]. The local routing is proposed in [10] to modify the call setup process when the caller and callee are in one VLR (Visitor Location Register), so the cost between the originating MSC (Mobile Switch Center) and the terminating HLR (Home Location Register) can be saved. The work in [15] uses a local routing policy for call setup based on three-tier database architecture in 3G network with the caching in GLR (Gateway Location Register).

Furthermore, the researches in IMS try to improve session setup delay in two aspects. The first is optimizing Session Initiation Protocol (SIP) signaling transmitting. IMS session setup delay is affected by the quality of the wireless link, e.g. frame error rate (FER), which can result in retransmissions of lost packets and can lengthen the session setup time. One way to do is choosing the appropriate retransmission timer and the underlying protocols. The work in [16] focuses on SIP signaling transmitting by optimizing it with an adaptive retransmission timer and evaluates SIP session setup performances with various underlying protocols, such as transport control protocol (TCP), user datagram protocol (UDP), and radio link protocols (RLPs). The work in [17] proposes that choosing an appropriate SIP compression efficiency and transport protocols can improve session setup delay. The work in [18] studies the SIP signaling transmitting, processing and queuing delay in 3G and WiMax networks, and proposes increasing channel rates can reduce IMS session setup delay.

The second is improving SIP signaling flows. The work in [8] investigates the call control procedure in UMTS Packet Switch (PS), and decreases call setup delay through performing Radio Access Network (RAN) resource allocation concurrent with media negotiation. But signaling interactions are reduced at the high cost of air interface resource to achieve fast handshake. The work in [19] is concerned about the in-calling setup delay and enhances the I-CSCF reliability through check point mechanism; also it uses the cache in I-CSCF to accelerate session setup. To the problem of triangular routing for a certain period of time when the user is

moving, Alam, M. T. *et al.* [20] proposes a decision algorithm to select the optimal session setup option.

The features of call service in IMS are similar to those in traditional 2G and 3G network, such as characteristics of localization. However, the previous works seldom take advantage of the locality to improve IMS session setup. Therefore, we propose the cache based session setup mechanism involving the local routing policy along with an adaptive caching algorithm, in order to solve the problem of signaling waste and reduce session setup delay.

### 3. Session Setup in IMS

#### 3.1. Basic Session Setup Flow

Figure 1 shows a classic session setup procedure in IMS network. With no loss of general we could suppose that as follows:

1) UE A and UE B are IMS terminals with the same type of properties.

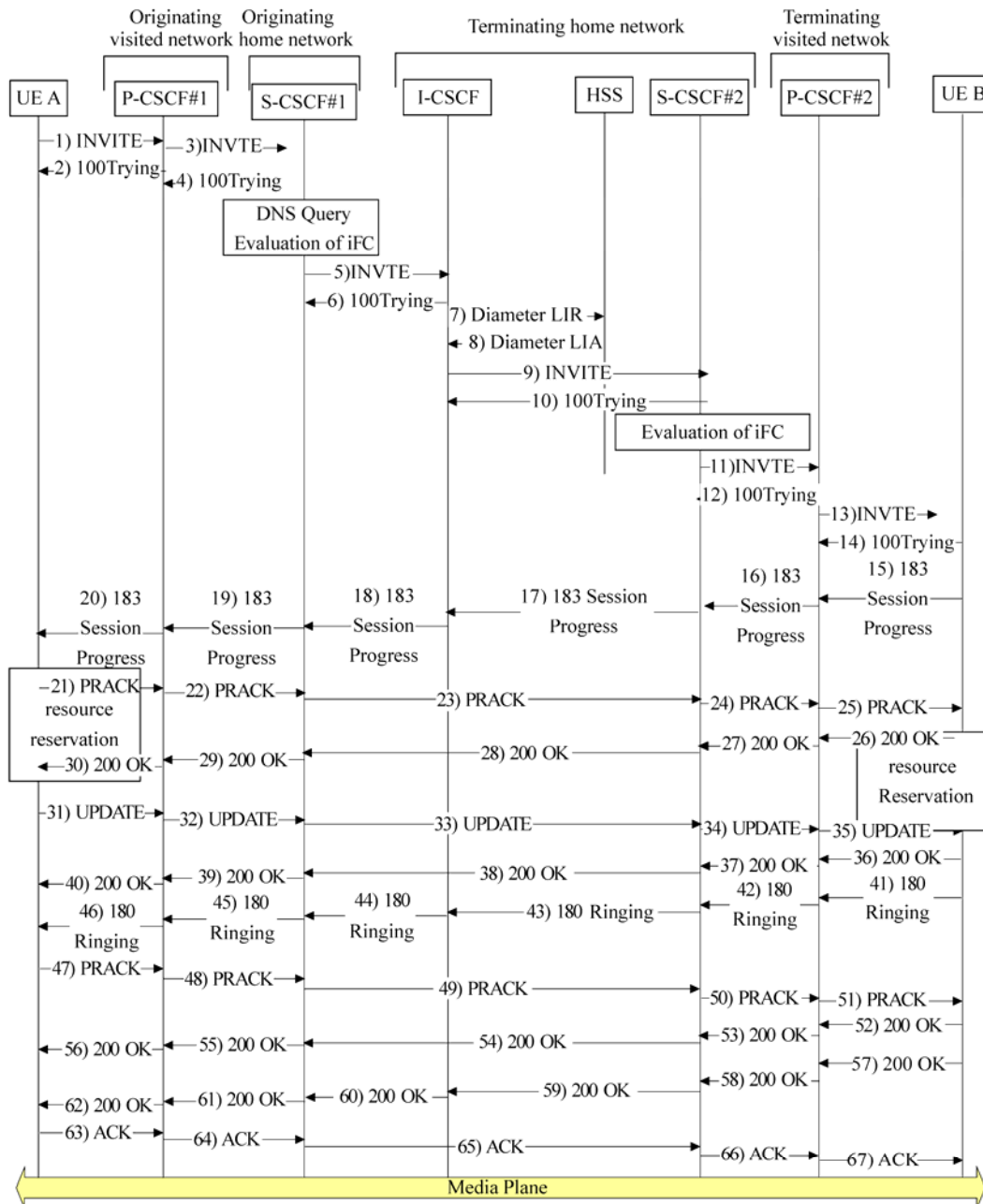


Figure 1. Session establishment procedure in IMS.

2) For simple consideration, both the caller and callee have no service contacts with the session.

3) P-CSCF # 1 and S-CSCF # 1 in the originating network are the entities of P-CSCF and S-CSCF providing services to the caller. Similarly, P-CSCF # 2 and S-CSCF # 2 in terminating network are the entities of P-CSCF and S-CSCF providing services to the callee.

As shown in Figure 1, when originating user A wants to establish session with terminating user B, UE A initiates a call to UE B by sending SIP 'INVITE'. As the 'INVITE' request arrives at S-CSCF # 1 via P-CSCF # 1, S-CSCF # 1 controls the services and the session. This procedure of each flow is given below: verify initial filter criteria (iFC); determine the address of network entrance the I-CSCF by querying DNS; forward 'INVITE' request to I-CSCF; the I-CSCF queries HSS to obtain the address of S-CSCF # 2, and then forwards the 'INVITE' request to S-CSCF # 2; finally when S-CSCF # 2 completes verifying iFC, the 'INVITE' request is forwarded to terminating UE B via P-CSCF # 2. Originating and terminating users take Quality of Service (QoS) negotiation by Session Description Protocol (SDP), which is carried in SIP message. After the two sides' QoS negotiation, resource reservation (3-40 steps in Figure 1) is processed. When UE A completes resource reservation, it notifies UE B by 'UPDATE' request, and UE B responses '200 OK' to confirm it. Following UE B's completion of resources reservation, '180 Ringing' message is sent back to the originating UE for ringing. When the terminating user answers, UE B sends the '200 OK' which is the response of 'INVITE' request. As soon as UE A receives this response, 'ACK' is sent to confirm (41-67 steps in Figure 1) that the session between originating and terminating users is established.

### 3.2. Determine the Address of the Terminating S-CSCF

In IMS basic session setup procedure, the originating S-CSCF (S-CSCF#1) is the first node that tries to forward the SIP request based on destination address which is in

the 'Request-URI' field carried by SIP 'INVITE'. The P-CSCF and the I-CSCF are not concerned the destination address, which means they don't inspect 'Request-URI' field in the SIP request. So the originating S-CSCF is the first point parsing the destination address, that is, according to the 'Request-URI' of SIP request it determines the next-hop address in the terminating network. During this procedure, the originating S-CSCF may find two different types of 'Request-URI': 'SIP URI' or 'TEL URI'. If the 'SIP URI' is found, a normal SIP process is adopted and 'INVITE' request is forwarded to the I-CSCF in terminating network through multi-steps DNS querying to determine next SIP server's address, which consists of transport protocol, hostname and port number that the I-CSCF supports [21]. And if the 'TEL URI' is found in 'Request-URI', DNS ENUM (E.164 number and DNS) is needed to decide the right I-CSCF address in terminating network. After receiving the 'INVITE' request, the I-CSCF gets the terminating S-CSCF (S-CSCF#2) address from the HSS by Diameter LIR (Location-Information-Request) and Diameter LIA (Location-Information-Answer) messages and relays 'INVITE' request. Thus, as shown in Figure 2 (A), the route between originating and terminating network has been set up. The detailed signalling flow presented in the specification [6] for the origination S-CSCF towards the terminating S-CSCF is summarized as follows. 1) If the analysis of the destination address determined that it belongs to a subscriber of a different operator, the request is forwarded to a well-known entry point in the destination operator's network, i.e. the I-CSCF. Then the I-CSCF queries the HSS for current location information and forwards the request to the S-CSCF. 2) If the analysis of the destination address determines that it belongs to a subscriber of the same operator, the S-CSCF forwards the request to a local I-CSCF, who queries the HSS for current location information. Then, the I-CSCF forwards the request to the S-CSCF.

Obviously, in originating network so many DNS queries and terminating S-CSCF discovery processes have seriously impact on the session setup delay and the cost of network transport. For IMS network which serves tens of thousands of users, performance will be significantly

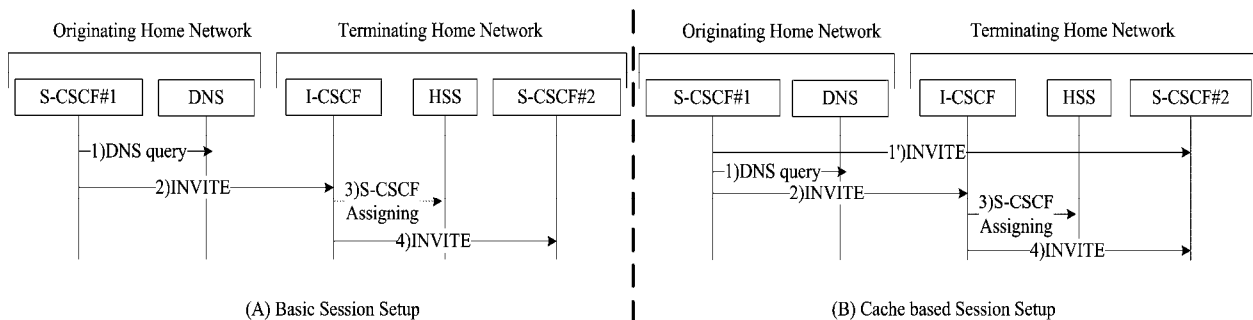


Figure 2. To find the terminating S-CSCF.

improved if we try to save a few messages for per session. Therefore, optimizing of IMS basic session setup mechanism is took into account. Furthermore, by addressing the terminating S-CSCF from the originating S-CSCF fleetly we can decrease the session setup delay and save network resources.

## 4. The Cache Based Session Setup Mechanism

This part describes the details of cache based session setup mechanism. The principles of this mechanism are: 1) If caller and callee are located in one S-CSCF serving area, their session will be established directly; 2) the originating S-CSCF forwards the 'INVITE' requests directly to the terminating S-CSCF, instead of the I-CSCF through DNS query; 3) the originating S-CSCF caches the terminating S-CSCF address information obtained from the previous session, and set a valid time for the cached information, which is an exponential distribution.

### 4.1. Cache Based Local Routing Policy

According to the principles of the cache based session setup mechanism, the local routing policy is adopted during session setup. The originating S-CSCF first checks whether the callee UE has registered on the originating S-CSCF or not. If true, session is established directly. If originating and terminating users are not in one S-CSCF, the originating S-CSCF checks whether the current terminating S-CSCF addresses has been cached or not, then it can determine that the next-step is directly visiting the terminating S-CSCF or querying DNS for retrieving the I-CSCF addresses.

The details are as follows, see Figure 2 (B):

1) When the originating S-CSCF receives the session setup request from caller, it checks whether the callee is within the same S-CSCF currently (the UE learns which the S-CSCF will be serving it through the IMS registration). If originating and terminating users are in one S-CSCF, session is established directly; else goes to 2.

2) The originating S-CSCF queries local cache to look for whether there has been the address of the S-CSCF serving the callee in previous sessions. If not, the S-CSCF performs a DNS query to retrieve the I-CSCF address, and sends 'INVITE' request to the terminating S-CSCF via the I-CSCF for session establishment, otherwise goes to 3. The originating S-CSCF caches the terminating S-CSCF address from the first response message (e.g. '183' message). 1–4 steps in Figure 2 (B).

3) For the originating S-CSCF cached the terminating S-CSCF address information, the originating S-CSCF sends 'INVITE' request to the terminating S-CSCF to

establish session, 1' step in Figure 2(B)

4) When the terminating S-CSCF receives originating 'INVITE' request, re-registration process perhaps be initiated because of the terminating user's roaming. A new S-CSCF is selected again in terminating home network, and the S-CSCF (the current terminating S-CSCF) is no longer available with returning error message. The originating S-CSCF needs to re-initiate basic session setup process. 1–4 steps in Figure 2 (B).

### 4.2. Adaptive Caching Algorithm

How to cache this information decides the accuracy of address query for the S-CSCF serving the terminating user. When the session setup request arrives, if cache information is effective, the originating S-CSCF hit the terminating S-CSCF with the most prefect performance. If the cached information in terminating S-CSCF is invalid or inexistence, basic session establishment procedure is needed, then at least the performance of cache based session setup is not worse than that of the basic one. But if cached information is outdated, after session arriving at a wrong location, it should re-establish in accordance with basic procedure, which is obviously the worst. Therefore, the effective of cached information has a great impact on probability of hitting terminating location in session establishment procedure.

Consequently, we design an adaptive caching algorithm in terms of the mobility patterns and the locality of call traffic rules, i.e. the probability that caller and callee locate in the same S-CSCF serving area is large. And we assume there is a data buffers in S-CSCF and the buffer size can satisfy system requirement.

Let  $G$  be the set of all out-area users.

$A = \{x|x : \text{the out-area users with new location in-information}\}.$

$B = G - A = \{y|y : \text{the out-area users with comparative mobile stabilization}\}.$

Define vector  $V$ . //  $V$  expresses the state of out-area user address information.

If  $V(x) = \text{TRUE}$

$x \in B$ ;

Else

$x \in A$ .

End If

Define vector  $T(x)$  for the caching time of the in-information. //  $T(X)$  expresses the period from previous resetting to the present time.

Assume  $Z$  is the out-area user obtained in the session setup procedure.

If  $z \in G$

  Set  $G = G \setminus \{z\}$ ,  $B = B \cup \{z\}$ ,  $V(z) = \text{TRUE}$ ,  $T(z) = 0$ ;

Else



Compare the address information of this procedure and the previous cached one.

If they are same &  $V(z)=FALSE$

Move the address information of user  $z$  from set  $A$  to set  $B$ ;

Set  $V(z)=TRUE$ ,  $T(x)=0$ ;

Else If they are not same &  $V(z)=TRUE$

Move the address information of user  $z$  from set  $B$  to set  $A$

Set  $V(z)=FALSE$ .

End If

End If

The judgment of cached information is:

Let the cache time threshold be  $D$ .

Let out-area callee as  $w$ .

If  $w \in G$

If  $V(w)=TRUE$

If  $T(w) < D$

Get the address information of user  $w$  from set  $B$ .

Else If  $T(w) > D$

Address as basic session setup procedure

End If

Else If  $V(w)=FALSE$

Session is established according to basic method;

End If

Else If  $w \in G$

Session is established according to basic method.

End If

## 5. Analytical Model

This section deduces the cache based session setup cost function and mean-delay function. For simplicity, we

have only analyzed the process shown in Figure 2, without considering the delay, the cost of UE, P-CSCF entities etc. Our definition of parameters is shown in Table 1.

According to [14],  $\alpha, \beta, \gamma$  are defined.  $\alpha$  means the probability of cache valid, i.e. the address of S-CSCF serving the terminating user has been cached in the originating S-CSCF and the called UE is in the terminating S-CSCF service area as it receives the session request.  $\beta$  means the probability of cache invalid: the originating S-CSCF hasn't cached the address information of S-CSCF serving the callee.  $\gamma$  means the probability of cache miss, i.e. the address of the terminating S-CSCF serving callee has been cached, but the callee no longer resides the terminating S-CSCF. According to the definition, we know  $\alpha + \beta + \gamma = 1$ . The cost function of basic session setup mechanism is given by

$$C_{basic} = C_{dns} + C_{s-i} + C_i + C_{i-s} + C_{s-cscf} \quad (1)$$

The mean delay function of basic session setup mechanism is:

$$D_{basic} = d_{dns} + d_{s-i} + d_i + d_{i-s} + d_{s-cscf} \quad (2)$$

The cost function of cache based session setup mechanism is:

$$C_{caching} = P_l \times C_{s-cscf} + (1 - P_l) \times (\alpha \times (C_{s-s} + C_{s-cscf}) + \beta \times (C_{dns} + C_{s-i} + C_i + C_{i-s} + C_{s-cscf}) + \gamma \times (C_{s-s} + C_{dns} + C_{s-i} + C_i + C_{i-s} + C_{s-cscf})) \quad (3)$$

The mean delay function of cache based session setup mechanism is:

$$D_{caching} = P_l \times d_{s-cscf} + (1 - P_l) \times (\alpha \times (d_{s-s} + d_{s-cscf}) + \beta \times (d_{dns} + d_{s-i} + d_i + d_{i-s} + d_{s-cscf}) + \gamma \times (d_{s-s} + d_{dns} + d_{s-i} + d_i + d_{i-s} + d_{s-cscf})) \quad (4)$$

**Table 1. Parameters definition.**

<i>Symbol</i>	<i>Quantity</i>	<i>Value</i>
$C_{dns} / d_{dns}$	The cost / mean time delay for performing a DNS querying	15u/15t
$C_{s-i} / d_{s-i}$	The cost / mean time delay for transmitting message from the originating S-CSCF to the terminating I-CSCF	10u/10t
$C_i / d_i$	The cost / mean time delay for one process of the S-CSCF assignment by the I-CSCF	20u/20t
$C_{i-s} / d_{i-s}$	The cost / mean time delay for transmitting message from the I-CSCF to the S-CSCF in the terminating network	5u/5t
$C_{s-s} / d_{s-s}$	The cost / mean time delay for transmitting message from the originating S-CSCF to the terminating S-CSCF	25u/25t
$C_{s-cscf} / d_{s-cscf}$	The cost / mean time delay of the S-CSCF	30u/30t
$P_l$	The probability of caller and callee in one S-CSCF	
$\alpha$	The probability of cache hit	
$\beta$	The probability of cache invalid	
$\gamma$	The probability of cache miss	

Then we try to derive probabilities of cache valid, miss, and invalid to (3), (4). Let the user residence time in S-CSCF as  $t_s$  which is assumed as an exponential distribution with parameter  $\lambda_s$ , and its probability density function is:

$$f_s(t) = \lambda_s e^{-\lambda_s t} \quad (5)$$

Denote by  $t_c$  the interval between two consecutive calls to the terminator, and  $t_m$  the interval between the arrival of previous call and the time as terminator move out the S-CSCF service area.  $f_c(t)$  and  $f_m(t)$  are density functions of  $t_c$  and  $t_m$ . We assume that the incoming call is a Poisson process, and then we have

$$f_c(t) = \lambda_c e^{-\lambda_c t} \quad (6)$$

According to the random observer property, we have

$$f_m(t) = \lambda_s \int_{r=t}^{\infty} f_s(r) dr = \lambda_s e^{-\lambda_s t} \quad (7)$$

We assume that the resident time of cached data is an exponential distribution with parameter  $\lambda_h$ , so

$$f_h(t) = \lambda_h e^{-\lambda_h t} \quad (8)$$

While caller initiates a session, if the address of S-CSCF serving the callee has been cached in the originating S-CSCF and the callee is still in cached S-CSCF, then the result is cache valid, and  $\alpha$  is

$$\begin{aligned} \alpha &= P[t_c < t_m \cap t_c < t_h] \\ &= \int_{t_m=0}^{\infty} \int_{t_c=0}^{\infty} \int_{t_h=0}^{\infty} f_m(t_m) f_c(t_c) f_h(t_h) dt_h dt_c dt_m = \frac{\lambda_c}{\lambda_s + \lambda_h + \lambda_c} \end{aligned} \quad (9)$$

If the originating S-CSCF has already removed the cached information of the address of S-CSCF serving the callee before an incoming call, the result is cache invalid  $\beta$ . The probability  $\beta$  is given by

$$\beta = 1 - P_r[t_c < t_h] = 1 - \int_{t_h=0}^{\infty} \int_{t_c=0}^{t_h} f_h(t_h) f_c(t_c) dt_c dt_h = \frac{\lambda_h}{\lambda_c + \lambda_h} \quad (10)$$

As  $\alpha + \beta + \gamma = 1$ , we can get

$$\gamma = 1 - \alpha - \beta = \frac{\lambda_h + \lambda_s}{\lambda_c + \lambda_h + \lambda_s} - \frac{\lambda_h}{\lambda_c + \lambda_h} \quad (11)$$

## 6. Performance Analysis

In this section, we firstly verify the validity of analytical model by using simulation experiments, and then we use numerical examples to investigate the performance of the proposed cache based session setup mechanism. For calculation convenience, let basic unit of cost be  $u$ , and basic unit of time delay be  $t$ . Parameter values are shown in Table 1. As the session setup cost of Equation (3) are the same as the session setup delay of Equation (4) in form, we only take the simulation experiments for the metric of session setup cost.

### 6.1. Verify Analytical Results with Simulation Results

In our simulation, there are provided the IMS network topology consisting UEs, the P-CSCF, the I-CSCF, the S-CSCF and the HSS. The simulation signaling flows are the same as Figure 2. The situations of user roaming and session initiating are simulated by generating discrete-events, including three types: 1) call event; 2) cache update; and 3) UE roaming. To investigate the impact of various network parameters on the performance of the new mechanism, the probability of roaming and the probability of caller and callee in one S-CSCF service area are varied by using different simulation configurations.

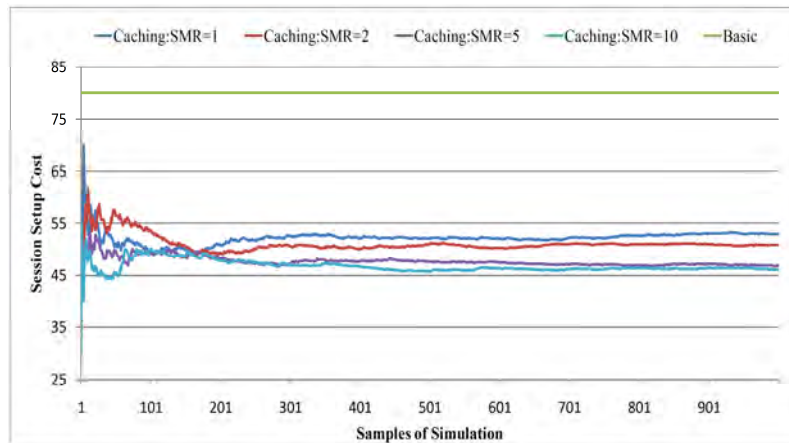


Figure 3. Simulation results of session setup cost.

**Table 2. Simulation and analytical results with difference of cost ( $P_l=0.5$ ).**

$\lambda_c$ (1/s)	$SMR(\lambda_c/\lambda_s)$	$C_{caching}$		Error (%)
		Simulation	Analytical	
800	1	54.9999	55	0.00%
400	2	53.81361761	54.053	0.44%
266.66667	3	50.7577782	50.739	-0.04%
200	4	49.20909	48.862	-0.71%
160	5	47.78306	47.674	-0.23%
133.33333	6	47.1684	46.858	-0.66%
114.28571	7	46.3097	46.263	-0.10%
100	8	46.5301	45.811	-1.57%
88.88889	9	45.771	45.456	-0.69%
80	10	45.3963	45.169	-0.50%

Here, the  $\lambda_s$  is five time of  $\lambda_h$ , and the SMR (Session to Mobility Ratio), expressed as  $\lambda_c/\lambda_s$ , is varied from 1 to 10. Then the cost values of two session setup mechanisms are measured and the values of 1000 samples got from the experiment in different network configuration are shown in Figure 3. Moreover, Table 2 shows the session setup cost values of the cache based session setup mechanism in both simulation experiments and numerical results, respectively. The cost values between simulation and model have some discrepancy due to the number of random generated discrete-events. And the jitter of the simulation values is also depicted in the Figure 3. If several more session events generated during the simulation period, the cost value of simulation is bigger than that of analytical value, and vice versa. Although the values between simulation and model have some discrepancy, as the error rates are all under 3%, these experiments have verified that analytical model is consistent with the simulation results.

## 6.2. Session Setup Cost and Delay

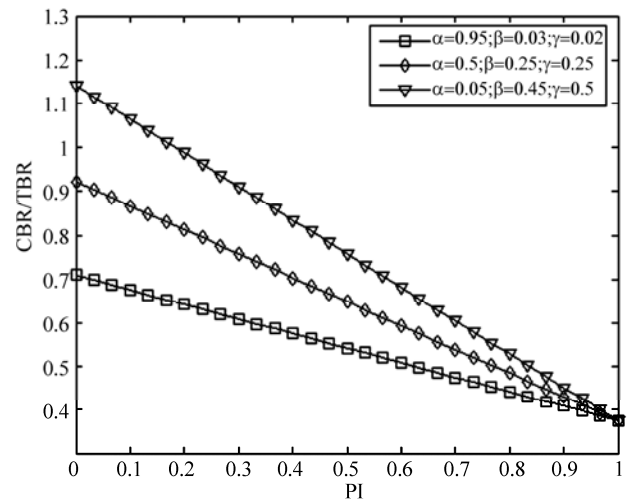
According to [14,15], we defined CBR (Cost Benefit Ratio) that is the ratio of  $C_{caching}$  to  $C_{basic}$  and TBR (Time Benefit Ratio) that is the ratio of  $D_{caching}$  to  $D_{basic}$ . We calculate session setup cost value and mean delay value in accordance with the formulas in Section 5. The conditions and results have been depicted in the figures.

Compared result curves are given in Figure 4 upon three different probabilities of  $\alpha, \beta, \gamma$ , with ratio results at y-axis and x-axis expressing the probability changing of the caller and callee in one S-CSCF. If  $P_l$  increases, the ratio of CBR and TBR becomes smaller. According to y-axis, along with the increasing of cache valid probability, the advantages of cache based session setup mechanism become more obvious. When the probability

of cache miss is bigger as  $\gamma=0.5$ , the performance of cache based session setup mechanism is worse than the performance of basic session setup mechanism ( $CBR/TBR>1$ ), because of session re-establishment after requests arriving the wrong location. However, we note that CBR/TBR has more than 1 just upon  $P_l < 0.18$ .

Considering the locality of call traffic (the number of local users calling local users is a large part total call number), the probability of callers and callees not in the same S-CSCF is smaller, so this situation will occur less. Therefore, generally speaking, cached based session setup mechanism outperforms the basic one.

According to the definition of SMR, the small value of SMR means the high mobility that UE has, vice versa. The curves varying along with the variable value of SMR and different  $P_l$  have been given in Figure 5. We can see that CBR and TBR decrease as SMR increases, which means caching contributes to decrease the session

**Figure 4. Comparison of CBR/TBR with caller and callee in one S-CSCF.**

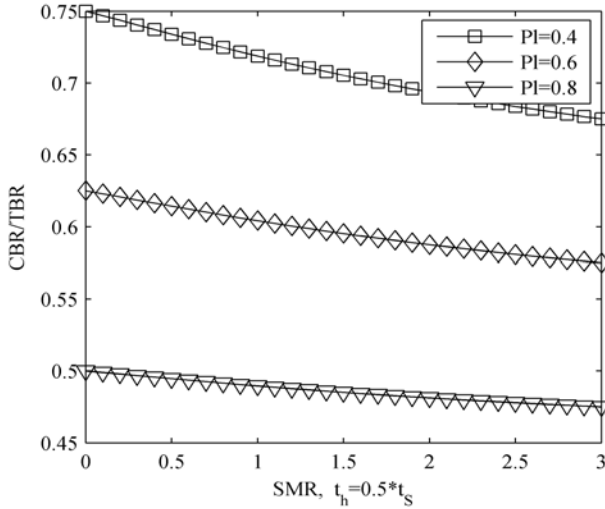


Figure 5. Comparison of CBR/TBR with  $t_h = 0.5t_s$ .

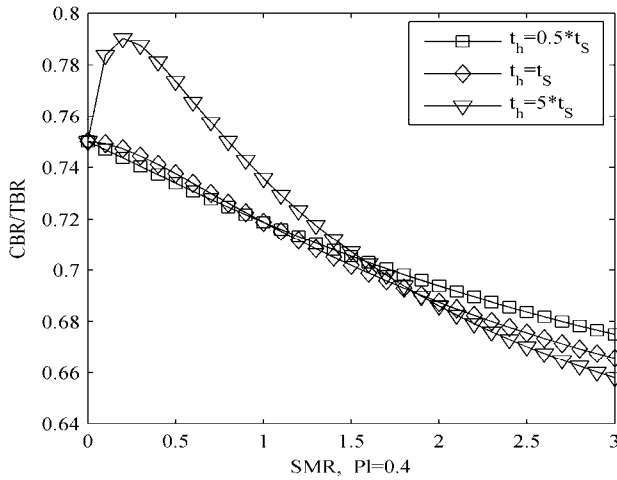


Figure 6. Comparison of CBR/TBR with  $P_l = 0.4$ .

delivery cost and mean delay if UEs have low mobility. Figure 6 presents the curves of CBR and TBR with three different multiple values of residence time and the cached data  $t_h$ . Study from y-axis, when the value of SMR is small (high mobility), the longer the caching time is, the greater the values of CBR and TBR are, and as the value of SMR is higher (low mobility) and the caching time is shorter, the values of CBR and TBR become the greater.

### 6.3. Query Accuracy Analysis

Let  $\alpha = 0, \beta = 0, \gamma = 1$ ,  $\text{CBR/TBR} = 1$ , then  $P_l = 0.333$ . In the worst situation  $n$ , i.e. cached information is outdated. After session arrives at wrong location, session

should be re-established in accordance with basic procedure. When  $P_l \leq 0.333$ , the ratio of CBR and TBR is equal to 1 or less than 1. This shows that enlargement of S-CSCF serving area conduces to improve the performance of session setup mechanism based on caching. We define the query accuracy as  $\eta = \alpha + \beta$ , then we get  $\eta + \gamma = \alpha + \beta + \gamma = 1$ . From the above analysis, the performance of session setup mechanism based on caching can be improved by increasing the query accuracy  $\eta$ . When caller and callee are in one S-CSCF serving area, no information required to cache and the performance of session setup procedure based on caching must be better than the basic. Then we discuss the required query accuracy when caller makes a call to the user in other S-CSCF serving area. For description convenience, we name the callee in this situation as out-area user.

From the definition of  $\eta$  and the Equations (9) and (10), we can derive the expression as follow:

$$\eta = \alpha + \beta = 1 - \gamma = \frac{\lambda_c}{\lambda_c + \lambda_h} \quad (12)$$

As shown in Figure 7, as SMR increasing,  $\eta$  increases. This is the query accuracy is added along with the value of CBR/TBR decreasing, i.e. the performance of session setup procedure based on caching is improved. Upon  $t_h = 0.5t_s$ , when SMR equals 0.333, the value of  $\eta$  can be got as 0.957, i.e. the probability of cache hit is large. And  $\eta$  decreases a little along with the SMR increasing. Upon  $t_h = 5t_s$ , when SMR equals to 0.333, the value of  $\eta$  is 0.592.  $\eta$  increases fast along with the SMR rising, and the query accuracy  $\eta$  is as large as 0.90 when SMR is 8.33. Thus, comparing of the two situations, to the out-area users with less mobility, the

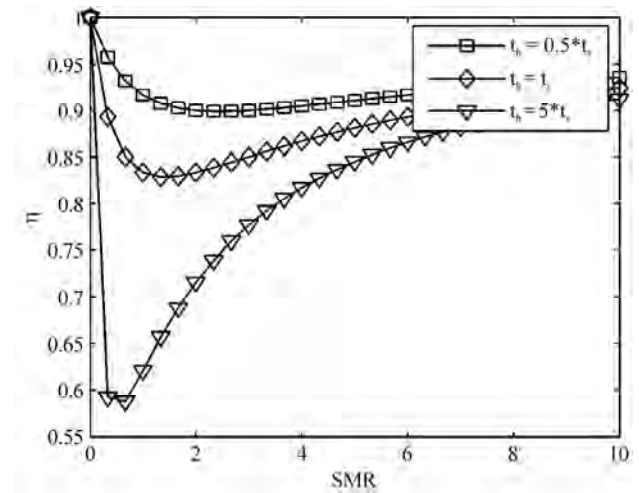


Figure 7. Query accuracy.

session setup mechanism based on cache has advantages of performance. The query accuracy will increase and the performance of cache based session establishment procedure will be improved while the effective time of cached data is prolonged.

## 7. Conclusions

How to decrease the network load and session delay is an important issue for designing and deploying IMS network. This article proposes: a cache based session setup mechanism, mainly through improving signaling flows in session establishment and introducing the caches to reduce network load and session delay. As it mostly perfects signaling flows, with no change of IMS core network architecture, no change or adding to terminal signaling, the cost of this improvement is comparatively smaller. Therefore, this innovation is practical for building of IMS network.

## 8. Acknowledgements

This work was jointly supported by: 1) National Science Fund for Distinguished Young Scholars (No. 60525110); 2) National 973 Program (No. 2007CB307100, 2007CB307103); 3) Development Fund Project for Electronic and Information Industry (Mobile Service and Application System Based on 3G).

## 9. References

- [1] H. Y. Dong, S. Kim, C. Nam, *et al.*, "Fixed and mobile service convergence and reconfiguration of telecommunications value chains," *Wireless Communications*, IEEE, Vol. 11, No. 5, pp. 42–47, October 2004.
- [2] 3GPP TS 23.228, V.8.2.0, IP Multimedia System (IMS), September 2007.
- [3] K. Shuang and F. C. Yang, "Study on enhanced 3G core network architecture," *ACTA ELECTRONICA SINICA*, Vol. 34, No. 7, pp. 1189–1193, July 2006.
- [4] 3GPP TS 22.228, V8.2.0. Service requirements for IP multimedia core network subsystem (IMS), October 2007.
- [5] 3GPP TS 23.218, V7.7.1. IP Multimedia(IM) session handling; IM call model, September 2007.
- [6] 3GPP TS 24.228, V5.15.0, Signaling flows for the IP multimedia call control based on SIP and SDP, October 2006.
- [7] 3GPP TS 24.229, V8.1.0. IP multimedia call control protocol based on SIP and SDP, September 2007.
- [8] K. Umschaden, I. Miladinovic, S. Bessler, and I. Gojmerac, "Performance optimizations in UMTS switched call control," *Fifth IEEE International Conference on 3G mobile communication technologies*, pp. 73–177, 2004.
- [9] G. Foster, M. I. Pous, D. Pesch, A. Sesmun, and V. Kenneally, "Performance estimation of efficient UMTS packet voice call control," *Proceedings of the IEEE Fall Vehicular Technology Conference*, pp. 1447–1451, 2002.
- [10] H. Jiang, B. Lu, and L. M. Li, "A novel call setup mechanism in a GSM network," *Journal of China Institute of Communications*, Vol. 23 No. 8, pp. 52–58, October 2002.
- [11] 3GPP TS 23.119, V7.0.0, Gateway Location Register (GLR), July 2007.
- [12] X. Jiang and I. F. Akyildiz, "A novel distributed dynamic location management scheme for minimizing signaling costs in mobile IP," *IEEE Transactions on Mobile Computing*, Vol. 1, No. 3, pp. 163–175, July–September 2002.
- [13] W. Ma and Y. Fang, "Improved distributed regional location management scheme for mobile IP," *IEEE Proceedings for Personal, Indoor and Mobile Radio Communications*, (PIMRC'03), Beijing, China, Vol. 3, pp. 2505–2509, September 7–10, 2003.
- [14] C. W. Pyo, J. Li, and H. Kameda, "A caching scheme for dynamic location management in PCS network," *IEEE 58th Vehicular Technology Conference*, Lbaraki, Japan, Vol. 2, pp. 761–765, 2003.
- [15] H. Zhang, J. X. Liao, and X. M. Zhu, "A novel call setup mechanism based on three-tier databases in 3G," *Journal of Electronics & Information Technology*, Vol. 29, No. 6, pp. 1290–1294, June 2006.
- [16] H. Fathi, S. S. Chakraborty, and R. Prasad, "Optimization of SIP session setup delay for VoIP in 3G wireless networks," *IEEE Transactions on Mobile Computing*, Vol. 5, No. 9, pp. 1121–1132, September 2006.
- [17] M. Melnyk and A. Jukan, "On signaling efficiency for call setup in all-IP wireless networks," *IEEE International Conference on Communications (ICC)*, Vol 5, pp. 1939–1945, June 2006.
- [18] A. Munir, "Analysis of SIP-based IMS session establishment signaling for WiMax-3G networks," *Fourth International Conference on Networking and Services (icns)*, pp. 282–287, 2008.
- [19] Y. B. LIN and M. H. TSAI, "Caching in I-CSCF of UMTS IP multimedia subsystem," *IEEE Transactions on Wireless Communications*, Vol. 5, No.1, pp. 186–192, January 2006.
- [20] M. T. Alam and Z. D. Wu, "Comparison of session establishment schemes over IMS in mobile environment," *2005 Fifth International Conference on Information, Communications and Signal Processing*, pp. 638–642, December 2005.
- [21] J. Rosenberg and H. Schulzrinne, *Session Initiation Protocol (SIP): Locating SIP Servers*, RFC 3263, IETF, June 2002.

# Metrics and Algorithms for Scheduling of Data Dissemination in Mesh Units Assisted Vehicular Networks\*

Zhongyi LIU, Bin LIU, Wei YAN

School of EECS, Peking University, Beijing, China

E-mail: {lzy, liubin, yanwei}@net.pku.edu.cn

Received March 19, 2009; revised May 8, 2009; accepted May 12, 2009

## Abstract

Data dissemination is an important application in vehicular networks. We observe that messages in vehicular networks are usually subject to both time and space constraints, and therefore should be disseminated during a specified duration and within a specific coverage. Since vehicles are moving in and out of a region, dissemination of a message should be repeated to achieve reliability. However, the reliable dissemination for some messages might be at the cost of unreliable or even no chance of dissemination for other messages, which raises tradeoffs between reliability and fairness. In this paper, we study the scheduling of data dissemination in vehicular networks with mesh infrastructure. Firstly, we propose performance metrics for both reliability and fairness. Factors on both the time and space dimensions are incorporated in the reliability metric and the fairness in both network-wide and Mesh Roadside Unit-wise (MRU-wise) senses are considered in the fairness metric. Secondly, we propose several scheduling algorithms: one reliability-oriented algorithm, one fairness-oriented algorithm and three hybrid schemes. Finally, we perform extensive evaluation work to quantitatively analyze different scheduling algorithms. Our evaluation results show that 1) hybrid schemes outperform reliability-oriented and fairness-oriented algorithms in the sense of overall efficiency and 2) different algorithms have quite different characteristics on reliability and fairness.

**Keywords:** Data Dissemination, Vehicular Networks, Scheduling, Mesh Backhaul

## 1. Introduction

Recently, vehicular networks with the assistance of roadside units (RSUs) have received considerable attention [1–5]. RSUs are useful in many different scenarios, such as Internet access “on the go”, collecting of sensed data from the sensors on vehicles, buffering data at hotspots, etc. However, we propose vehicular networks with mesh backhaul. As wireless mesh networks have the potential advantage of easy deployment, self-configurability and large coverage [6], mesh routers are adequate to act as RSUs, which we call MRUs (Mesh Roadside Units).

In vehicular networks, data are often subject to some type of time constraints and space constraints. For example, congestion information is meaningless for vehicles 10 miles away and might become invalid after two hours.

Other types of messages can include the following: “Road maintenance work will be performed from 3:00pm to 4:00pm at the Lincoln Street”, “Traffic control will be enforced from 10:00am to 10:30am near the railway station”. Messages can also be generated by the transportation monitoring system, such as “The Lincoln Street is often in congestion from 5:00pm to 6:00pm”. This kind of messages should be disseminated during a specified *duration* and within a specific *coverage*. As vehicles are constantly moving in and out of a region, dissemination of a message should be repeated in the specified duration to achieve reliability, namely to ensure that at all times all the vehicles in a region are notified. However, the reliable dissemination of some messages might be at the cost of the unreliable or even no chance of dissemination of other ones, which raises tradeoffs between *reliability* and *fairness*. In this scenario, reliability has the meaning in both time and space dimensions. The time dimension depicts the reliability achieved by a specific MRU in its scheduling process while the space

\*This work is co-supported by National Key Basic Research Program of China (No. 2009CB320504 and No. 2007CB310902), State key lab. of virtual reality technology and systems and Peking University-Morgan Stanley Research Fund.

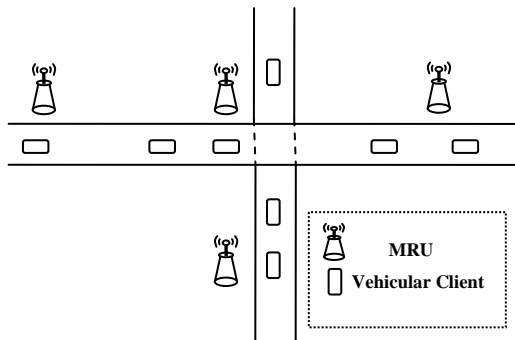
dimension describes whether messages are disseminated at all the MRUs within their requested coverage. Similarly, fairness is significant in both the network-wise sense and the MRU-wise sense. In this paper, we first propose metrics for reliability and fairness and then develop and evaluate five scheduling algorithms quantitatively. To the best of our knowledge, this is the first paper on scheduling mechanisms in this scenario to address the reliability and fairness issues.

The rest of this paper is arranged as follows: the second part presents our system model, including the network architecture and performance metrics. Five different scheduling algorithms are proposed in Section 3 and evaluation results are shown in Section 4. Related work is reviewed in Section 5, and in the last section, we conclude this paper.

## 2. System Model

### 2.1. Network Architecture

We assume a vehicular network with mesh backhaul, as shown in Figure 1. Mesh Roadside Units (MRUs) are placed at roadside locations to receive messages from cars nearby or to disseminate information to vehicles in its vicinity. Since wireless mesh network has the potential advantage of easy deployment and self-configurable, we assume that MRUs are connected via wireless links. MRUs can have larger coverage than vehicular clients, so that they can serve more vehicular clients at the same time and disseminate messages efficiently. When a vehicle needs to send messages to an MRU, it can first select a nearby MRU and then looks for relays to forward information to the designated MRU. However, the mechanism with which vehicles send messages to MRUs is out of the scope of this paper. We assume perfect message transmission from vehicular clients to MRUs in this work. We focus on the scheduling for data dissemination in this type of vehicular networks, which will be explained in detail in later sections.



**Figure 1. Network architecture.** We assume a mesh backhaul assisted architecture. Mesh road side units are assumed to be well connected.

### 2.2. Performance Metrics

In our message dissemination model, each message is coupled with a  $\langle \text{start-time, end-time, } x, y, \text{radius} \rangle$  tuple, in which “start-time” and “end-time” indicate the instants when message dissemination should begin and terminate;  $\langle x, y \rangle$  and “radius” specify the center and radius of the dissemination area. With “x”, “y”, “radius” and some geographical information, the MRU which receives the message dissemination request can easily obtain the destination MRUs which locate in the destination area. With this message dissemination model, we figure out two important factors which determine the efficiency of message dissemination:

- **Reliability.** Reliability describes the quality of service for the selected messages. In this case, reliability covers two different dimensions. On one hand, in the time dimension, a message should be given as much dissemination time as possible in the [start-time, end-time] duration. On the other hand, in the space dimension, message dissemination should occur in an area as large as possible within the requested  $\langle x, y, \text{radius} \rangle$  coverage.
- **Fairness.** To achieve better reliability for the a selected message to disseminate, more time as well as MRUs should be allocated to it, which may in turn decrease the quality of service for the other messages. Therefore, fairness should be taken into account to enhance the efficiency of message dissemination.

We now present the metrics for both reliability and fairness. Our metric for reliability, RM (reliability metric) is defined as

$$RM = \alpha * TRM + (1 - \alpha) * SRM \quad (0 \leq \alpha \leq 1)$$

where TRM and SRM stand for Time Reliability Metric and Space Reliability Metric, respectively. The formulas for TRM and SRM are as follows.

$$TRM = \frac{\sum_{m \in MRU_i} TR(m, i)^\kappa}{N_{MRU}^i}$$

$$SRM = \frac{\sum SR(m)^\kappa}{N_{messages}}$$

where  $N_{MRU}$  is the total number of MRUs in the network and  $N_{MRU}^i$  is the number of messages to disseminate (or received) at the  $i$ th MRU.  $TR(m, i)$  indicates the *time ratio* for message  $m$  at the  $i$ th MRU and  $SR(m)$  is the *space ratio* for message  $m$ . The exponent  $\kappa \geq 1$  is used to specify the reliability level, whose effects will be explained later in this section.  $TR(m, i)$  and  $SR(m)$

are in turn defined as  $TR(m, i) = \frac{D\_time(m, i)}{Duration(m)}$  and



$SR(m) = \frac{N_{D\_MRU}^m}{N_{C\_MRU}^m}$ , respectively, namely  $TR(m, i)$  is the ratio between dissemination time allocated to message  $m$  at the  $i$ th MRU and the requested duration of message  $m$ , and  $SR(m)$  is the ratio between number of MRUs which provide dissemination service for  $m$  and the overall number of MRUs within the requested  $\langle x, y, \text{radius} \rangle$  coverage. It is clear that our metric of reliability incorporates the reliability factors in both the time dimension and the space dimension.

The selection of  $\kappa$  is critical to achieve different reliability levels. Take the definition of TRM as an example. We assume the dissemination cycle of an MRU is  $\Delta$ , which we call it a *slot*. Assume there are two messages  $m_1, m_2$  to disseminate, both with the same duration of 2 slots and the same coverage. If  $\kappa = 1$ , then the scheduling sequences  $[m_1, m_2]$  (the first slot for  $m_1$  and the second slot for  $m_2$ ) and  $[m_1, m_1]$  will result in the same TRM. This is because in the first sequence

$$TR(m_1) + TR(m_2) = \frac{1}{2} + \frac{1}{2} = 1$$

And in the second sequence

$$TR(m_1) + TR(m_2) = 1 + 0 = 1$$

However, if we choose  $\kappa = 2$ , then the TRM of the two sequences will be different. For the first one,

$$TR(m_1) + TR(m_2) = \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 = \frac{1}{2}$$

And for the second one

$$TR(m_1) + TR(m_2) = (1)^2 + (0)^2 = 1$$

Therefore, the sequence with *less* messages will achieve higher reliability. The value of  $\kappa$  for TRM and SRM can be different. However, in this paper, we assume the same reliability level is used for both TRM and SRM.

Our metric for fairness, FM (fairness metric) is defined as

$$FM = \beta * \frac{N_D}{N_R} + (1 - \beta) * \frac{\sum_{MRU_i} \frac{N_D^i}{N_R^i}}{N_{MRU}} \quad (0 \leq \beta \leq 1)$$

In which  $N_D$  and  $N_R$  stand for the number of disseminated messages and the number of dissemination requests for the network while  $N_D^i$  and  $N_R^i$  stand for the number of disseminated messages and the number of dissemination requests at the  $i$ th MRU. It should be clear that our fairness metric combines fairness factors in both the network-wise sense and the MRU-wise sense. However, currently our fairness metric only reflects whether a message is given the opportunity to be disseminated, without regarding whether different messages are given

the same level of opportunities. We leave this topic as our future work.

Also, we can combine the two metrics together. Therefore the combined metric, CM, can be defined as

$$CM = \gamma * RM + (1 - \gamma) * FM \quad (0 \leq \gamma \leq 1)$$

### 3. Scheduling Algorithms

Given the system model described above, we developed several scheduling algorithms, which exhibit different characteristics of reliability and fairness. As been stated before, we assume the dissemination cycle of an MRU is  $\Delta$  and call it a slot. A given duration between start-time and end-time can be transformed into the equivalent representation with the ordinal number of dissemination cycles. The task of scheduling algorithms is to determine the message to disseminate in the future  $W$  dissemination cycles, here we call  $W$  the *schedule window*. We further assume that any MRU knows locations of all the MRUs in the network, so that a given geographical coverage can be mapped into an equivalent representation with a list of MRUs. In the following sections, the coverage of a message means the number of MRUs in the geographical area. All the scheduling algorithms have a time complexity of  $O(W * n)$ , where  $W$  is the size of schedule window and  $n$  is the number of messages to disseminate.

#### 3.1. MQIF-Maximum Quality Increment First

Our first scheduling algorithm, Maximum Quality Increment First (MQIF) scheduling, serves first the messages which would bring the maximum quality of service (namely reliability) increment. Our approach is to first calculate the expected increment in TRM and then estimate the expected increment in SRM assuming allocation the current cycle to a message. Afterwards, the two increments are combined. The detail of MQIF is shown in Figure 2.

For each slot in the future schedule window, MQIF compares the expected quality increments of all the messages whose duration covers that slot and selects the one with the maximum quality increment. The precise calculation of expected increment of RM (denoted as QI) assuming the allocation of a slot to a message is impossible at runtime in a distributed manner. Therefore, we use the scheme shown in Figure 3 to estimate the value of QI.

The expected quality increment (QI) can be obtained from the expected increment of TRM (denoted as TQI) and that of SRM (denoted as SQI). TQI can be easily got from local information. However, SQI is dependent on the scheduling results of other MRUs in the coverage of the given message. Since we don't assume one MRU knows the scheduling status of other MRUs, we estimate

```

MQIF_Schedule()
1. selected_messages ← []
2. for i ← 1 to schedule_window do
3.   max_QI ← 0
4.   for j ← 1 to number_messages do
5.     if (msg[j].start_time - current_slot ≤ i)
        AND (msg[j].end_time - current_slot ≥ i) then
6.       QI ← calc_QI(msg[j])
7.       if QI > max_QI then
8.         max_QI = QI
9.         selected_messages[i]=msg[j]
10.      end if
11.    end if
12.  end for
13.  if selected_messages[i] ≠ null then
14.    selected_messages[i].selection_count ++
15.  end if
16. end for

```

Figure 2. Maximum quality increment first (MQIF) scheduling.

```

Calc_QI(msg)
1.  $TQI \leftarrow \left(\frac{msg.selection\_count+1}{msg.duration}\right)^\kappa - \left(\frac{msg.selection\_count}{msg.duration}\right)^\kappa$ 
2.  $r \leftarrow random(0, msg.coverage)$ 
3. if msg.selection_count = 0 then //estimate increment of SRM
4.    $SQI \leftarrow \left(\frac{r+1}{msg.coverage}\right)^\kappa - \left(\frac{r}{msg.coverage}\right)^\kappa$ 
5. else
6.    $SQI \leftarrow 0$ 
7. end if
8.  $QI \leftarrow \alpha * TQI + (1-\alpha) * SQI$ 
9. return QI

```

Figure 3. Calculating the expected increment of reliability metric.

the current number of MRUs who have already scheduled the given message by generating a random number in the range 0 to *msg.coverage*.

```

LSF_Schedule()
1. selected_messages ← []
2. for i ← 1 to schedule_window do
3.   min_selection_count ← INFINITY
4.   for j ← 1 to number_messages do
5.     if msg[j].start_time - current_slot ≤ i
        AND msg[j].end_time - current_slot ≥ i then
6.       selection_count ← msg[j].selection_count
7.       if selection_count < min_selection_count then
8.         min_selection_count = selection_count
9.         selected_messages[i]=msg[j]
10.      end if
11.    end if
12.  end for
13.  if selected_messages[i] ≠ null then
14.    selected_messages[i].selection_count ++
15.  end if
16. end for

```

Figure 4. Least selected first (LSF) scheduling.

### 3.2. LSF-Least Selected First

The second scheduling approach, Least Selected First (LSF) scheduling, tries to schedule into the schedule window as many messages as possible. The general idea is that if a message had the least opportunity to be served before, it will be given the highest priority this time. LSF is given in Figure 4.

For each slot in the future schedule window, LSF compares the selection count of all the messages whose duration covers that slot and allocate the slot to the message with the minimum selection count.

### 3.3. Hybrid Schemes

Since MQIF tends to achieve high reliability and LSF tends to achieve good fairness, we can combine the two strategies to make tradeoffs between the two metrics. We figure out two approaches to do this:

- Add a certain condition to MQIF or LSF. We call the resulting algorithm Conditional-MQIF or Conditional-LSF. In Conditional-MQIF, MQIF strategy is applied only when the given condition is met; otherwise the LSF strategy is adopted. Different conditions can result in different tradeoffs between reliability and fairness; therefore this approach can be adapted to different application scenarios easily.

Combine the two strategies by simply combining the selection metrics of the two. Although it is less tunable than the former one, it may achieve better overall efficiency.

```

Conditional_MQIF_Schedule()
1. selected_messages ← []
2. for i ← 1 to schedule_window do
3.   max_QI ← 0, min_QI ← INFINITY
4.   min_selection_count ← INFINITY
5.   for j ← 1 to number_messages do
6.     if msg[j].start_time - current_slot ≤ i
        AND msg[j].end_time - current_slot ≥ i then
7.       QI ← calc_QI(msg[j])
8.       selection_count ← msg[j].selection_count
9.       if QI > max_QI then
10.        max_QI = QI
11.        MQIF_msg ← msg[j]
12.      end if
13.      if QI < min_QI then
14.        min_QI = QI
15.      end if
16.      if selection_count < min_selection_count then
17.        min_selection_count = selection_count
18.        LSF_msg ← msg[j]
19.      end if
20.    end if
21.  end for
22.  //determine which strategy to use
23.  if  $\frac{\max\_QI}{\min\_QI} \geq MQIF\_THRESHOLD$  then
24.    selected_messages[i] = MQIF_msg
25.  else
26.    selected_messages[i] = LSF_msg
27.  end if
28.  if selected_messages[i] ≠ null then
29.    selected_messages[i].selection_count ++
30.  end if
31. end for

```

Figure 5. Conditional-MQIF.

### 3.3.1. Conditional-MQIF

The first approach to combine MQIF and LSF is to add a threshold to MQIF or LSF. In Conditional-MQIF, MQIF strategy is applied only when the ratio between the maximum and minimum QI exceeds a specified *MQIF\_THRESHOLD*. If the condition is not met, LSF is applied. Similarly, we can also combine MQIF and LSF using Conditional-LSF. In Conditional-LSF, LSF is used only when the difference between the maximum and minimum selection count exceeds a predefined *LSF\_THRESHOLD*. We show the detail of conditional-MQIF in Figure 5. By varying the *MQIF\_THRESHOLD* or *LSF\_THRESHOLD*, we can control the proportion of opportunities for applying MQIF or LSF strategy; therefore the algorithm can be adapted to various application demands. For example, small *MQIF\_THRESHOLD* values tend to achieve better reliability thus adequate for reliability-sensitive scenarios.

### 3.3.2. MQILSF-Maximum Quality Increment Least Selected First

Since MQIF tends to select messages with small coverage and duration, while LSF favors messages with small selection count, we can simply incorporate selection count into the quality increment (QI). This strategy is similar to MQIF, except that in MQILSF the Quality Increment (QI) is redefined as

$$QI = \frac{\alpha * TQI + (1 - \alpha) * SQI}{msg[j].selection\_count + 1}$$

Note that we use  $msg[j].selection\_count + 1$  instead of  $msg[j].selection\_count$  because the initial values of selection counts are all 0.

## 4. Performance Evaluation

We developed a discrete event simulator in Java. It takes as input an XML configuration file and a scenario file, runs the designated scheduling algorithms and writes the scheduling results to trace files.

### 4.1. Simulation Setup

We extract a 2500m\*2500m network scenario from a realistic geographical map of the TianAnMen district of Beijing, whose e-map is shown in Figure 6 [13]. We assume MRUs are evenly distributed with a distance of 300m along the roads. Therefore over 50 MRUs are deployed. The communication range of an MRU is assumed to be 350m.

The *message generation interval* and *message generation probability* indicate how often events are generated.

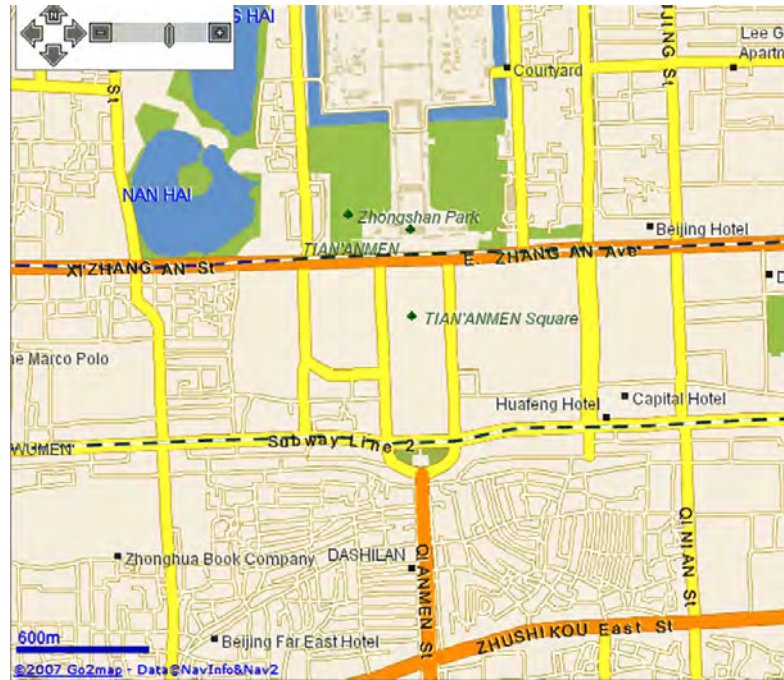


Figure 6. Simulated scenario.

Table 1. Simulation setup.

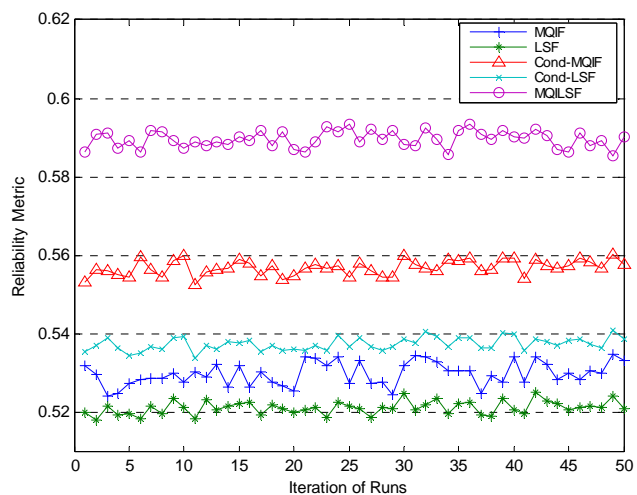
Parameter	Value
Simulation time	6000s
Dissemination Cycle of MRUs	1s
Duration of messages	5s~300s
Coverage of messages	600m~1500m
Message Generation Interval	30s
Message Generation Probability	0.15
Schedule Interval	5s
Reliability Level $\kappa$	2
Reliability Metric Parameter $\alpha$	0.5
Fairness Metric Parameter $\beta$	0.5
Combined Metric Parameter $\gamma$	0.5
Threshold in Conditional-MQIF	15
Threshold in Conditional-LSF	15

The simulation parameters are shown in Table 1. In our experiments, an opportunity is given to an MRU every 30 seconds and the MRU generates a message with a probability of 0.15. Durations of messages are in the range [5s, 300s] and the coverage of messages are in the range [600m, 1500m], namely about 2~5 hops. The reliability level  $\kappa$  is set to 2 in Subsection 4.2 and Subsection 4.4. Threshold values for Conditional-MQIF and Conditional-LSF are all set to 15 in Subsection 4.2 and 4.3.

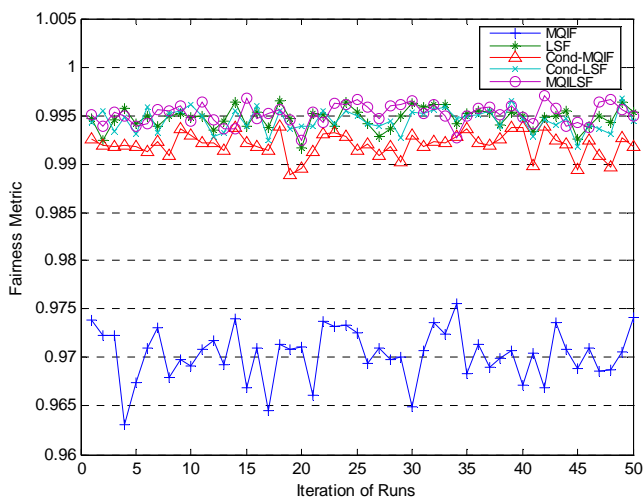
## 4.2. Comparison of Different Schemes

The Reliability Metric (RM), Fairness Metric (FM) and Combined Metric (CM) achieved by different scheduling algorithms are shown in Figure 7(a), Figure 7(b) and Figure 7(c), respectively.

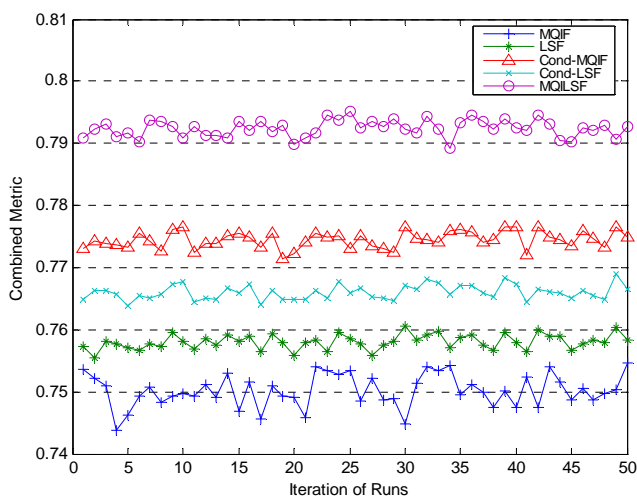
It is not hard to understand that the reliability metric of LSF and the fairness metric of MQIF are the worst among all the algorithms. However, it is interesting that



(a)

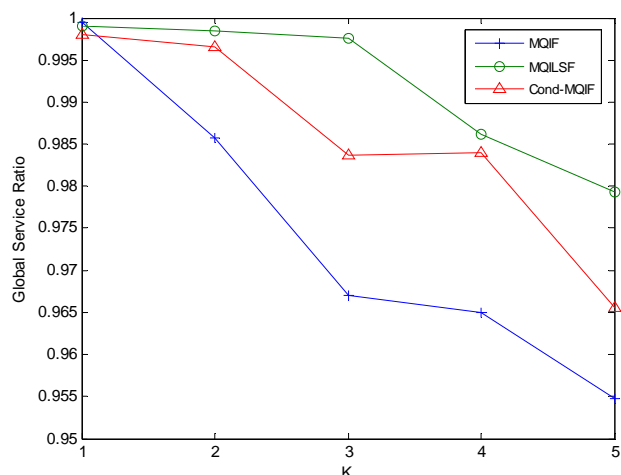


(b)

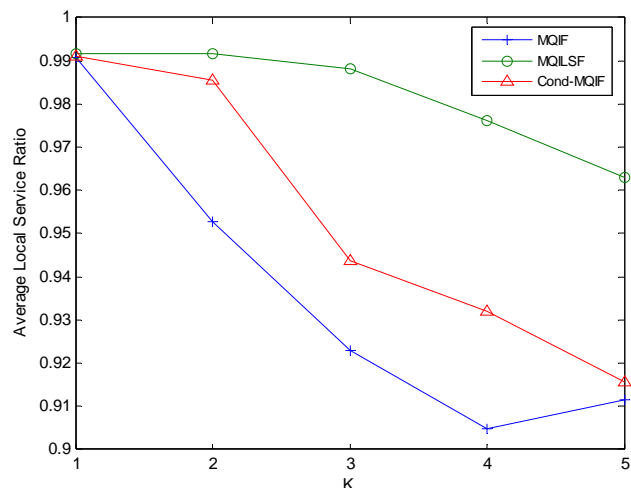


(c)

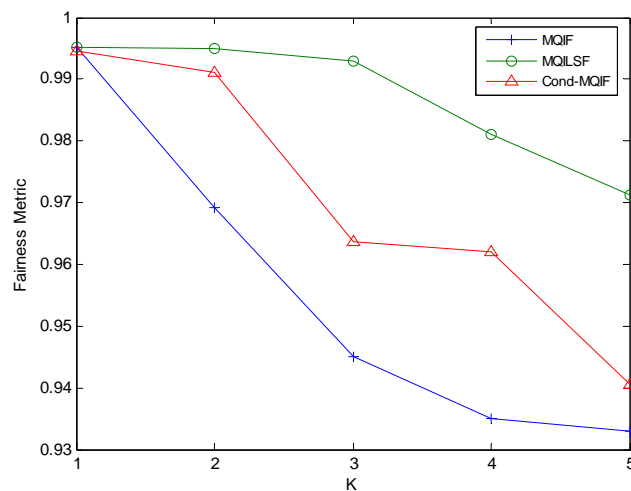
**Figure 7.** a) Comparison of reliability metric, b) Comparison of fairness metric, c) Comparison of combined metric.



(a)



(b)



(c)

**Figure 8.** a) Effects of  $\kappa$  on global service ratio, b) Effects of  $\kappa$  on average local service ratio, c) Effects of  $\kappa$  on fairness metric.

MQIF does not achieve the best reliability. We attribute this to the fact that MQIF is a greedy approach but its decisions are not made based on deterministic information. Hybrid algorithms achieve better reliability and fairness, therefore better combined metric. For example, the Reliability Metric of Conditional-LSF is about 7% higher than that of LSF and they achieve about the same level of fairness metric; the Reliability Metric and the Combined Metric of MQILSF are about 11% and 7.5% higher than those of MQIF, respectively.

### 4.3. Effects of $\kappa$

The different values of  $\kappa$  indicate different reliability levels. Although we cannot analyze the effects of  $\kappa$  on reliability by directly comparing the values the reliability metric under different models, we can do analysis by comparing the number of messages disseminated globally and locally. The Global Service Ratio (GSR), which is defined as

$$\text{GSR} = \frac{N_D}{N_R},$$

reflects the service ratio of the overall network, where  $N_D$  and  $N_R$  stand for the number of message dissemination requests received and the number of disseminated messages of the network, respectively. Note that GSR is actually the first part of the fairness metric. The Average Local Service Ratio (Average-LSR), which is defined as

$$\text{Average-LSR} = \frac{\sum_{MRU_i} \frac{N_D^i}{N_R^i}}{N_{MRU}},$$

reflects the average of the local service ratio of all MRUs in the network. Note that Average-LSR is actually the second part of the fairness metric. We also study the net effect of  $\kappa$  on the fairness metric.

As shown in Figures 8(a), (b) and (c), in MQIF, Conditional-MQIF and MQILSF, the Global Service Ratio, the Average Local Service Ratio and the Fairness Metric all decrease as  $\kappa$  increases.

Specially, the effect of  $\kappa$  on the Average Local Service Ratio is stronger than on the Global Service Ratio and the Fairness Metric, and moreover, MQIF is extremely sensitive to the value of  $\kappa$  while MQILSF is the least sensitive. This may indicate that MQILSF has the advantage of increasing reliability without degrading fairness very much.

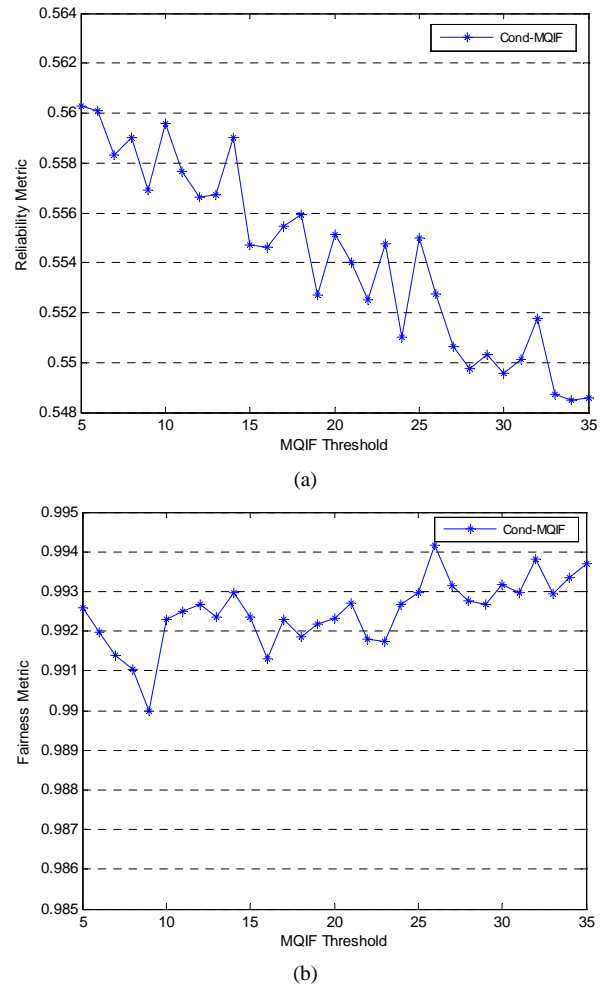
### 4.4. Effects of Threshold Values

In Conditional-MQIF and Conditional-LSF, the threshold values are critical on the reliability and fairness metric achieved.

As shown in Figure 9(a) and Figure 9(b), in Conditional-MQIF, the reliability metric decreases as the threshold value increases while the fairness metric increases as the threshold value increases. This can be attributed to the fact that the larger the threshold, the less opportunities MQIF strategy is adopted while the more opportunities are given to the LSF strategy. Similar trends are also observed in Conditional-LSF. In Conditional-LSF, as the threshold increases, more opportunities are given to the MQIF strategy, which results in better reliability metric and smaller fairness metrics. In our simulated scenario, the best threshold for Conditional-MQIF is 14 or 15, while the best for Conditional-LSF is any number in [15,18].

## 5. Related Work

Although a lot of work has been done to develop vehicular networks with infrastructure [1–5], they are usually restricted to one-hop communication between vehicular clients and roadside units. However, we propose using



**Figure 9. a) Effects of threshold values on reliability metric, b) Effects of threshold values on fairness metric.**

wireless mesh routes as the backhaul of the network, which has the potential advantage of easy deployment, self-configurable and scalability.

Scheduling for data access in vehicular networks is studied in [5]. However, our work is different from [5] because

- The work by [5] only studies scheduling for data access within one hop. In contrast, we focus on the multi-hop case, which is realistic for data dissemination in vehicular networks.
- The work by [5] does optimization for scheduling of upload/download data access. However, we consider the scheduling for data dissemination.
- The matter of fairness is not taken into account by [5]. We figure out that in data dissemination in our scenario, reliability and fairness should both be studied so that dissemination efficiency can be enhanced.

To the best of our knowledge, this is the first paper to address the reliability (in both the time dimension and the space dimension) and fairness issues in scheduling of data dissemination in the vehicular networks with mesh infrastructure.

A large amount of work has been performed on packet scheduling of MAC layer in wireless networks. The work by [7,8] tried to providing packet-level quality of service by packet scheduling. The main goals of [7,8] is to achieve fairness and maximum channel utilization. The work by [9] proposed OSMA, a packet scheduling approach in MAC layer to enhance throughput by choosing a receiver with good channel condition. However, none of them address the time and space constraints in the scenario of data dissemination in vehicular networks.

## 6. Conclusions and Future Work

As messages in vehicular networks are usually subject to space and time constraints, tradeoffs must be made between reliability and fairness for message dissemination algorithms. We propose the performance metrics for reliability and fairness in the scenario of scheduling for message dissemination in vehicular networks with mesh infrastructure. Five different scheduling algorithms are developed and evaluated quantitatively. We concluded that

- Although a greedy approach is adopted, the reliability-oriented algorithm, MQIF, does not achieve the best reliability. We attribute this to the fact that MQIF makes its greedy decisions based-on non-deterministic information.
- The fairness-oriented algorithm, LSF, achieves the best fairness metric as well as the worst reliability metric.
- The hybrid scheme, MQILSF, achieves the best reliability and combined metric and its fairness metric is nearly the same as LSF.

- The other two hybrid schemes, Conditional-MQIF and Conditional-LSF, are not as good as MQILSF. However, the idea of combining different algorithms by adding a certain condition to one algorithm can be helpful in other research fields, because it is easy to be adapted to different application scenarios.

Our evaluation on the reliability level parameter  $\kappa$  of the reliability metric show that different values of  $\kappa$  means different reliability levels. Therefore, for scenarios requiring different reliability levels, different values for  $\kappa$  should be used.

However, our current metric for fairness is not perfect; for example, it does not incorporate the relative dissemination time between different messages. On the other hand, different messages may have different priorities (indicating different level of importance or urgency), which is not considered in this paper. Furthermore, dynamic traffic densities may be useful for scheduling algorithms. For example, if the current traffic density is low, diversity of messages or fairness might be favored. Therefore, we plan to develop priority and traffic density aware scheduling schemes in the future.

## 7. References

- [1] V. Bychkovsky, B. Hull, *et al.*, "A measurement study of vehicular internet access using in situ wi-fi networks," In Proceedings of the 12th Annual International Conference on Mobile Computing and Networking (MOBICOM'06), pp. 50–61, 2006.
- [2] D. Hadaller, S. Keshav, T. brecht, *et al.*, "Vehicular opportunistic communication under the microscope," In Proceedings of the 5th International Conference on Mobile Systems, Applications, and Services (MobiSys'07), 2007.
- [3] B. Hull, V. Bychkovsky, Y. Zhang, *et al.*, "Cartel: A distributed mobile sensor computing system," In Proceedings of the 4th International Conference on Embedded Networked Sensor Systems (SenSys'06), pp. 125–138, 2006.
- [4] V. Navda, A. P. Subramanian, *et al.*, "MobiSteer: Using steerable beam directional antenna for vehicular network access," In Proceedings of the 5th International Conference on Mobile Systems, Applications, and Services (MobiSys'07), 2007.
- [5] Y. Zhang, J. Zhao, and G. H. Cao, "On scheduling vehicle-roadside data access," In Proceedings of the Fourth ACM International Workshop on Vehicular Ad Hoc Networks (VANET'07), pp. 9–18, 2007.
- [6] I. F. Akyildiz, X. Wang, *et al.*, "Wireless mesh networks: A survey," In Computer Networks, Vol. 47, No. 4, pp. 445–487, 2005.
- [7] H. Y. Luo, *et al.*, "A new model for packet scheduling in multihop wireless networks," In Proceedings of the 6th



- Annual International Conference on Mobile Computing and Networking (MOBICOM'00), pp. 76–86, 2000.
- [8] H. Y. Luo, *et al.*, “A packet scheduling approach to Qos support in multihop wireless networks,” In *Mobile Networks and Applications*, Vol. 4, pp. 193–206, 2004.
  - [9] J. F. Wang, *et al.*, “Opportunistic packet scheduling and media access control for wireless LANs and multi-hop ad hoc networks,” In *IEEE Wireless Communications and Networking Conference*, (WCNC'04), pp. 1234–1239, 2004.
  - [10] Y. Ding, *et al.*, “A static-node assisted adaptive routing protocol in vehicular networks,” In *Proceedings of the Fourth ACM International Workshop on Vehicular Ad Hoc Networks (VANET'07)*, pp. 59–68, 2007.
  - [11] I. Leontiadis, *et al.*, “Opportunistic spatio-temporal dissemination system for vehicular networks,” In *Proceedings of the 1st International Mobisys Workshop on Mobile Opportunistic Networking (MobiOpp'07)*, pp. 39–46, 2007.
  - [12] P. V. Kanodia, *et al.*, “Distributed multi-hop scheduling and medium access with delay and throughput constraints,” In *Proceedings of the 7th Annual International Conference on Mobile Computing and Networking (MOBICOM'01)*, pp. 200–209, 2001.
  - [13] E-map of Beijing (English Version), <http://en.beijing2008.cn/06/78/emap.shtml>.

# High Resolution MIMO-HFSWR Radar Using Sparse Frequency Waveforms

Guohua WANG, Yilong LU

*School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore*

*E-mail: {wang0330, eylu}@ntu.edu.sg*

*Received April 21, 2009; revised May 12, 2009; accepted May 15, 2009*

## Abstract

In high frequency surface wave radar (HFSWR) applications, range and azimuth resolutions are usually limited by the bandwidth of waveforms and the physical dimension of the radar aperture, respectively. In this paper, we propose a concept of multiple-input multiple-output (MIMO) HFSWR system with widely separated antennas transmitting and receiving sparse frequency waveforms. The proposed system can overcome the conventional limitation on resolutions and obtain high resolution capability through this new configuration. Ambiguity function (AF) is derived in detail to evaluate the basic resolution performance of this proposed system. The advantages of the system of fine resolution and low peak sidelobe level (PSL) are demonstrated by the AF analysis through numerical simulations. The impacts of Doppler effect and the geometry configuration are also studied.

**Keywords:** MIMO, HFSWR, Radar, Sparse Frequency Waveform

## 1. Introduction

HIGH frequency surface wave radar (HFSWR) is a low-cost radar system that adopts vertically polarized high frequency electromagnetic signals which propagate along the ocean surface. A preferable property of HFSWR is that it can detect and track ship and aircraft targets beyond the horizon. Due to this reason, HFSWR has a wide range of applications in both civil and military fields. For conventional HFSWR systems, its range resolution is highly restricted by the bandwidth of available clear channels in a congested spectrum environment [1,2], while the azimuth resolution is also constrained by the physical dimension of the radar antenna aperture.

Multiple-input multiple-output (MIMO) radar is now getting much intention for various applications such as detection, estimation, and imaging *etc.* MIMO radar can transmit at transmitters multiple waveforms that are dividual at the receivers so that it can obtain more degrees of freedom compared with conventional radars that transmit single waveform [3–6]. With widely distributed antennas, angular diversity can be fully achieved to compete with the target scintillations [7,8]. Meanwhile, MIMO radars with widely distributed antennas can gain high resolution by coherent processing [8]. Like distrib-

uted MIMO radar, a single-input multiple-output (SIMO) radar system with sparse coherent receiving aperture can also achieve high resolution on the order of one wavelength with limited bandwidth as reported in [9].

Sparse frequency waveform problem has been studied in [10,11] and literatures therein. For HFSWR, sparse frequency waveform can provide large flexibility to choose clear channels and thereby reduce interferences from assigned channels. Motivated by the potential benefits from sparse frequency waveforms and high resolution capacity of coherent multistatic radar and MIMO radar systems, we in this paper propose a novel MIMO-HFSWR using sparse frequency waveforms to break down the limitation on range and azimuth resolutions of conventional HFSWR. Ambiguity Function (AF) is derived in detail and fully investigated in this paper to analyze the performance of the proposed system. Unlike that of paper [8], we take Doppler effect in the AF for analysis. Through AF analysis it is demonstrated that this system has high flexibility in operation and attractive improvement on resolution in the restricted geographical condition as well as the congested spectrum environment. In particular, by using the widely separated antennas, it abates the aperture limitation as well as the rigorous land requirement successfully. In addition, by using sparse

frequency waveform it not only takes more clear channels into use to compete with co-channel interference but also reduces the peak sidelobe level (PSL).

The reminder of this paper is organized as follows. The AF of the proposed MIMO-HFSWR using sparse frequency waveforms is derived in Section 2. Based on AF analysis and simulations the system is evaluated in terms of resolution capacity and PSL performance, with zero Doppler frequency in Section 3 and under the factors of geometric configurations and Doppler effects in Section 4. Finally, conclusions and future work are outlined in Section 5.

## 2. Ambiguity Function of MIMO-HFSWR System Using Sparse Frequency Waveforms

Ambiguity function is an important tool in conventional radar analysis as it shows radar's inherent capacity of discriminating targets associated with different time delay and Doppler frequency. Thus, in this paper, we also employ AF to evaluate the performance of the proposed system.

The proposed MIMO-HFSWR system consists of  $M$  transmitters transmitting  $M$  waveforms and  $N$  receivers. Each transmitter is assigned a distinct channel with starting frequency  $f_m, m=1, 2, \dots, M$ . Thus, collectively, the transmitting waveforms will have a sparse spectrum, because which we call the transmitting waveforms sparse frequency waveform. All antennas are arbitrarily located with mutual separation distance larger than several wavelengths in a 3-dimensional space. Figure 1 shows the system configuration, where  $R_n$  and  $T_m$  refer to the  $n$ -th receiver and the  $m$ -th transmitter, respectively. Each transmitter and each receiver are located at a point represented by a 3-dimensional vector in the Cartesian coordinate system. For example, the  $m$ -th transmitter is associated with a vector  $\mathbf{c}_{t,m}=[x_m, y_m, z_m]$ , and the  $n$ -th receiver  $\mathbf{c}_{r,n}=[x_n, y_n, z_n]$ . For simplification, this paper considers only single point target case. And the target is assumed to be located at a general point  $\mathbf{x}=[x, y, z]$  with constant velocity of  $\mathbf{v}=[v_x, v_y, v_z]$ . As the antennas are widely distributed, each of them will view the target with

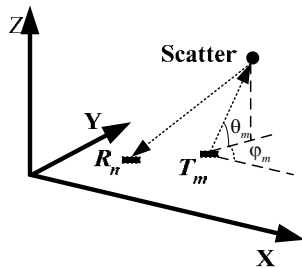


Figure 1. MIMO radar configuration.

a different angle. Angle variables  $\theta$  and  $\phi$  refer to the true elevation and azimuth as illustrated in Figure 1. We also assume that the phases and time at the transmitters and receivers are synchronized in advance. Meanwhile, the signal attenuation in different path is assumed to be the same.

Let  $x_m(t)$  be the signal transmitted by the  $m$ -th transmitter that meets the requirement of narrow band assumption. It is expressed as

$$x_m(t) = \exp(j2\pi f_m t) s_m(t) \quad (1)$$

where  $s_m(t)$  is the baseband waveform of the  $m$ -th transmitter. After the signal impinged back from the target to the  $n$ -th receiver, the echo is:

$$e_n(t) = \sum_{m=1}^M \gamma x_m(t - \tau_{nm}) \exp(-j2\pi f_{d_{nm}} t) \quad (2)$$

where  $\gamma$  is the complex reflection coefficient of the target,  $\tau_{nm}$  is the round-trip delay, and  $f_{d_{nm}}$  is the Doppler frequency of the echo at the  $n$ -th receiver due to the  $m$ -th transmitter. We take the assumption that the target stops during the pulse transmission and reception. Then  $\tau_{nm}$  is the round-trip delay at the start of observation time, and has the form

$$\begin{aligned} \tau_{nm} &= r_T^m(t)/c + r_R^n(t)/c \Big|_{t=0} \\ &= \sqrt{(x - x_m)^2 + (y - y_m)^2 + (z - z_m)^2}/c \\ &\quad + \sqrt{(x - x_n)^2 + (y - y_n)^2 + (z - z_n)^2}/c \end{aligned} \quad (3)$$

where  $c$  is the velocity of light in the media that the transmitters, receivers and targets are located in. In the case that the transmitters or the receivers are mounted on moving platform, the platform velocity can also be easily included in (3). For different scatter,  $\tau_{nm}$  is a function of variables  $x, y$ , and  $z$ . Thus, through Taylor-series analysis at a reference point  $\mathbf{x}_0=[x_0, y_0, z_0]$ , (3) can be changed to

$$\begin{aligned} \tau_{nm} &= r_T^m(0)/c + r_R^n(0)/c \Big|_{[x,y,z]} \\ &\quad \left( r_T^m(0)/c + r_R^n(0)/c \Big|_{[x_0,y_0,z_0]} \right) + \\ &\quad \left( \frac{(x - x_0)(\cos \theta_m \cos \phi_m + \cos \theta_n \cos \phi_n)}{c} + \right. \\ &\quad \left. \frac{(y - y_0)(\cos \theta_m \sin \phi_m + \cos \theta_n \sin \phi_n)}{c} + \right. \\ &\quad \left. \frac{(z - z_0)(\sin \theta_m + \sin \theta_n)}{c} \right) \\ &= \tau_{nm}^0 + \tau'_{nm} \end{aligned} \quad (4)$$

where

$$\tau_{nm}^0 = r_T^m(0)/c + r_R^n(0)/c \Big|_{[x_0,y_0,z_0]} \quad (5)$$

and

$$\tau'_{nm} = \left( \frac{(x-x_0)(\cos\theta_m \cos\varphi_m + \cos\theta_n \cos\varphi_n)}{c} + \frac{(y-y_0)(\cos\theta_m \sin\varphi_m + \cos\theta_n \sin\varphi_n)}{c} + \frac{(z-z_0)(\sin\theta_m + \sin\theta_n)}{c} \right) \quad (6)$$

$$f_{dnm} = \frac{v_x(x-x_m) + v_y(y-y_m) + v_z(z-z_m)}{\lambda_m \sqrt{(x-x_m)^2 + (y-y_m)^2 + (z-z_m)^2}} + \frac{v_x(x-x_n) + v_y(y-y_n) + v_z(z-z_n)}{\lambda_m \sqrt{(x-x_n)^2 + (y-y_n)^2 + (z-z_n)^2}} \\ = \frac{v_x(\cos\theta_m \cos\varphi_m + \cos\theta_n \cos\varphi_n)}{\lambda_m} + \frac{v_y(\cos\theta_m \sin\varphi_m + \cos\theta_n \sin\varphi_n)}{\lambda_m} + \frac{v_z(\sin\theta_m + \sin\theta_n)}{\lambda_m} \quad (8)$$

It can be easily proved that (8) is a tantamount expression of conventional Doppler frequency of bistatic radar.

As each transmitting waveform is assigned to a distinct channel, orthogonality holds for all the transmitting signals. Thus, at each receiver, signals from  $M$  transmit-

$$y_{nm}(\mathbf{x}, \mathbf{x}_0, \mathbf{v}) \\ = \gamma \exp(-j2\pi f_m(\tau_{nm} - \tau_{nm}^0)) \int_{-\infty}^{+\infty} s_m(t - \tau_{nm}^0) \exp(j2\pi f_{dnm}t) s_m^*(t - \tau_{nm}) dt + n_{nm}(t) \\ = \gamma \exp(-j2\pi f_m \tau'_{nm}) \mathbf{A}_{nm}(\tau'_{nm}, f_{dnm}) + n_{nm}(t) \quad (9)$$

where  $\mathbf{A}_{nm}$  is the correlation between the  $m$ -th transmitting waveform and its delay-Doppler shifted version, and  $n_{nm}(t)$  is the noise component of the output.

Collectively, there are  $NM$  outputs after matched filtering. By coherently summing all these outputs we can get:

$$\mathbf{A}(\mathbf{x}, \mathbf{x}_0, \mathbf{v}) = \left| \sum_{n=1}^N \sum_{m=1}^M y_{nm}(\mathbf{x}, \mathbf{x}_0, \mathbf{v}) \right|^2 \quad (10)$$

Because both  $\tau'_{nm}$  and  $f_{dnm}$  in (9) are affected by azimuth angles and elevation angels, the geometry configuration represented by a matrix  $\mathbf{C}$  consisting of all the azimuth and elevation angels should be included in the ambiguity function. Ignoring the noise-based component and discarding the target reflection coefficient in (10), we can define the normalized ambiguity function for the proposed system as:

$$\chi(\mathbf{x}, \mathbf{x}_0, \mathbf{v}, \mathbf{C}) = \frac{1}{M^2 N^2} \left| \sum_{n=1}^N \sum_{m=1}^M \exp(-j2\pi f_m \tau'_{nm}) \mathbf{A}_{nm}(\tau'_{nm}, f_{dnm}) \right|^2 \quad (11)$$

As  $\tau'_{nm}$  and  $f_{dnm}$  are related to  $x-x_0$ ,  $y-y_0$ ,  $z-z_0$ ,  $v_x$ ,  $v_y$ ,  $v_z$ , azimuth angles and elevation angels, the range and azimuth resolution as well as the effects of velocity components and geometry configuration can be assessed through the AF analysis.

$f_{dnm}$  has the form

$$f_{dnm} = \frac{1}{\lambda_m} \frac{d}{dt} (r_T^m(t) + r_R^n(t)) \quad (7)$$

Again, through Taylor-series analysis (7) can be changed to

ters can be firstly separated by down-converting into  $M$  channels. Then, for each channel, a matched filter of corresponding transmitting waveform is employed at the interested range cell centered at  $\mathbf{x}_0 = [x_0, y_0, z_0]$ . Thus, the  $m$ -th filter output at the  $n$ -th receiver can be expressed as

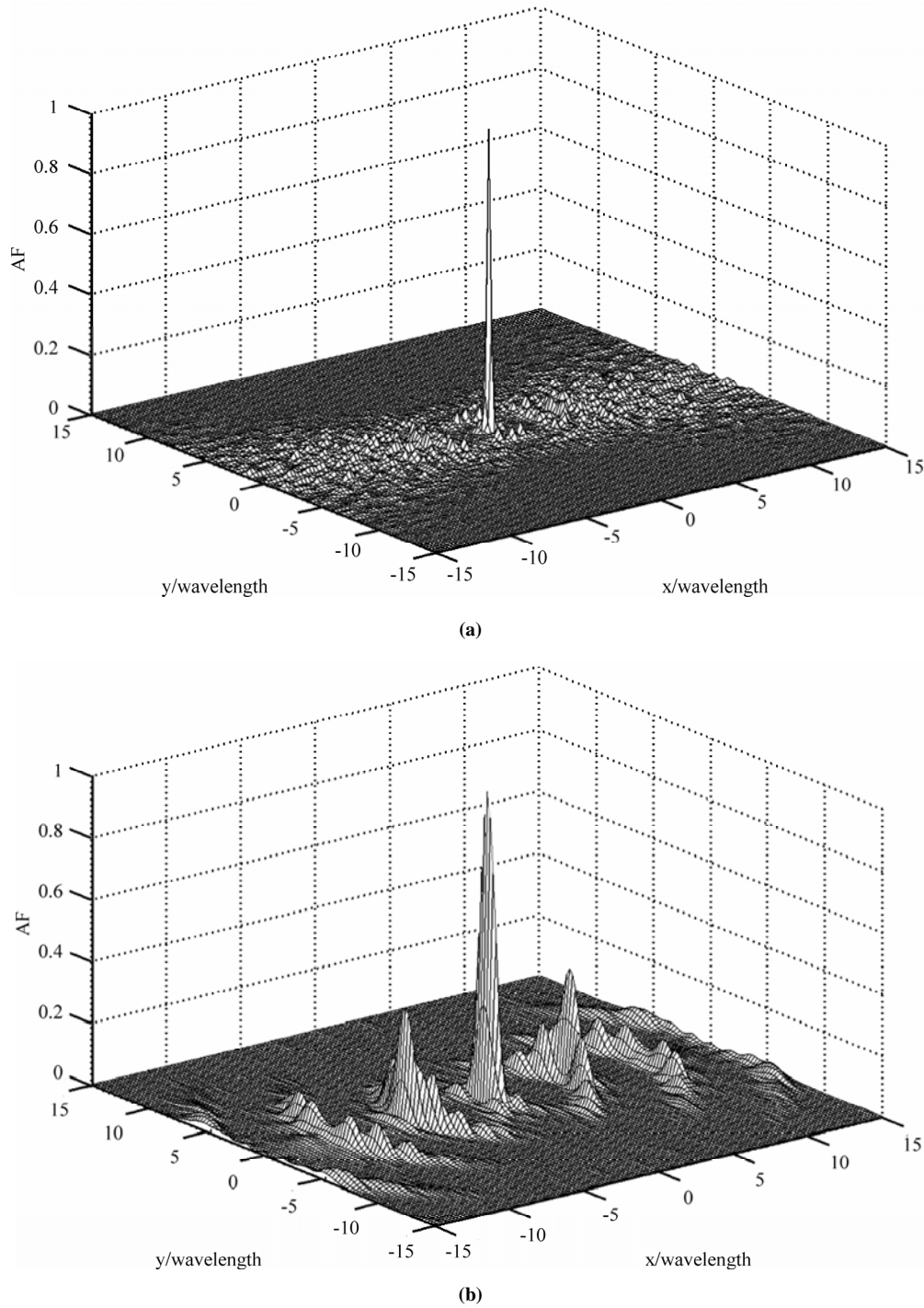
Besides, different waveforms may obtain different  $\mathbf{A}_{nm}$ , thus waveforms also play a key role in the MIMO radar ambiguity function. We take Linear Frequency Modulation (LFM) waveforms as an example to illustrate this point. A conventional LFM waveform defined by  $u(t) = \text{rect}(t/T) \exp(j\pi k t^2)$  has an correlation function like [12]:

$$\mathbf{A}(\tau, v)_{LFM} = \exp(-j\pi k \tau^2) \left( 1 - \frac{|\tau|}{T} \right) \sin c \left( (vT - B\tau) \left( 1 - \frac{|\tau|}{T} \right) \right) \quad (12)$$

where  $v$  is the Doppler frequency,  $\tau$  is the time-delay,  $k$  is the chirp rate,  $T$  is the pulse width,  $B$  is the bandwidth. From (12), it can be easily inferred that the bandwidth of single waveform adopted will impact the performance of the MIMO-HFSWR system.

### 3. Resolution Capacity and PSL Performance

In this section, the resolution capacity and PSL of the MIMO-HFSWR system using sparse frequency waveform are assessed by setting the Doppler frequency to zero in AF analysis. Sparse frequency waveforms consisting of stepped frequency linear frequency modulation signals are investigated. For simplification we just study a simplified 2-dimensional configuration.



**Figure 2. MIMO AF of (a) Sparse frequency waveforms (b) Single LFM signal, at zero Doppler frequency.**

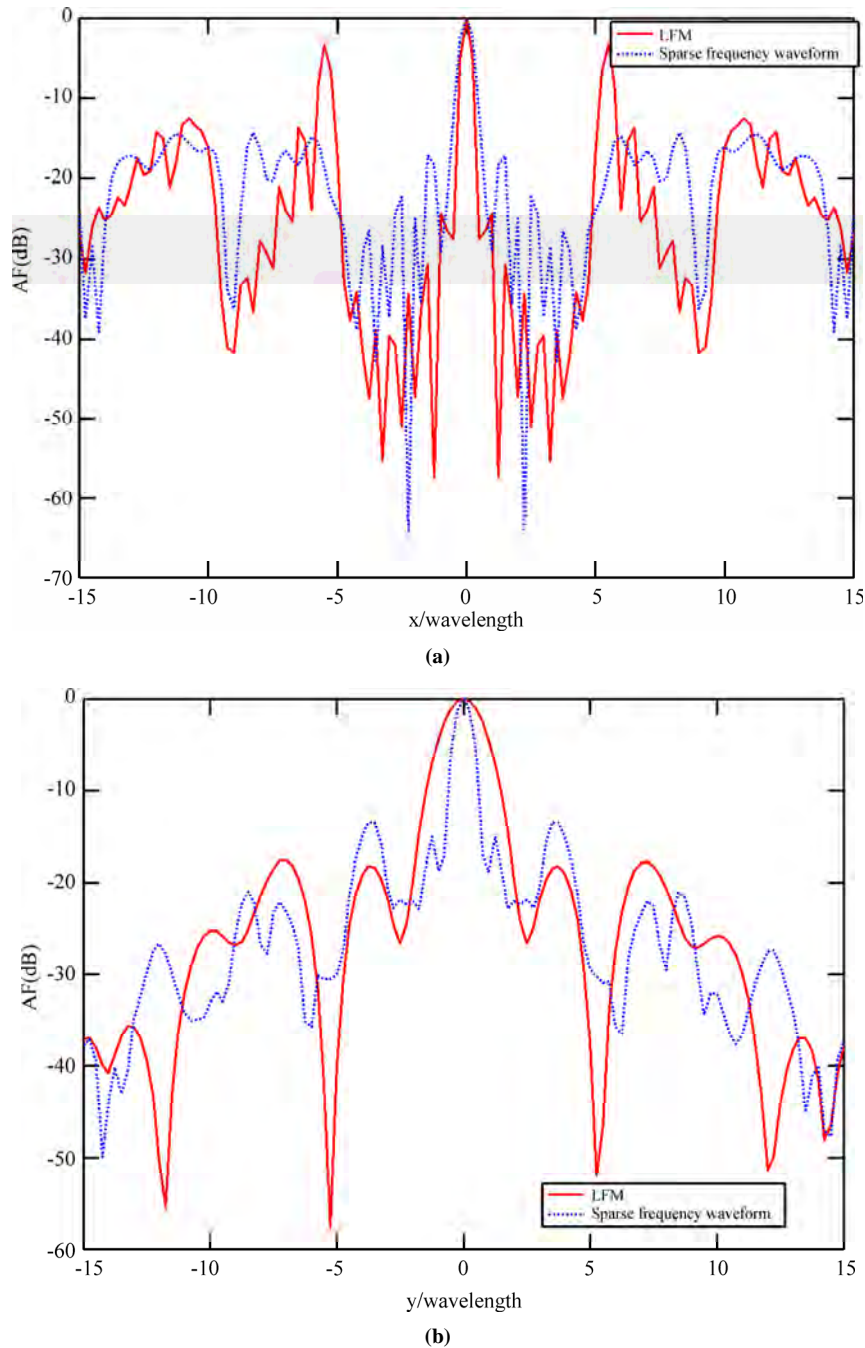
As we can see from above analysis, the ambiguity function of the proposed MIMO-HFSWR system depends on the system geometry configuration confined by all azimuth and elevation angles. Thus, we can ignore the true position of transmitters and receivers. We here take nine transmitters and nine receivers located evenly over

spatial region of  $(-\pi/4, \pi/4)$  for  $\varphi$ . Each transmitter will emit one LFM with an assigned start frequency. The pulse width is 100  $\mu$ s for all transmitters. The bandwidth of each LFM pulse is 500 kHz. The nine start frequencies of LFM waveforms are defined as the sequence of  $\{5, 6, 7, 8, 9, 8, 7, 6, 5\}$  MHz. Orthogonality can be achieved

by sequentially transmitting at transmitters or by setting the first five LFM waveforms to be up-chirps and the left four down-chirps [12]. We in this paper utilize the first mechanism. As a comparison, we also take an ambiguity function from a single LFM waveform with the same bandwidth and pulse width as well as the start frequency of 9 MHz. We also suppose to transmit it sequentially in time domain so that we can separate at each receiver the

returns from different transmitters. By central coherent processing we can also get the results of AF as showed in Figure 2(b), which seems the same as that in [8].

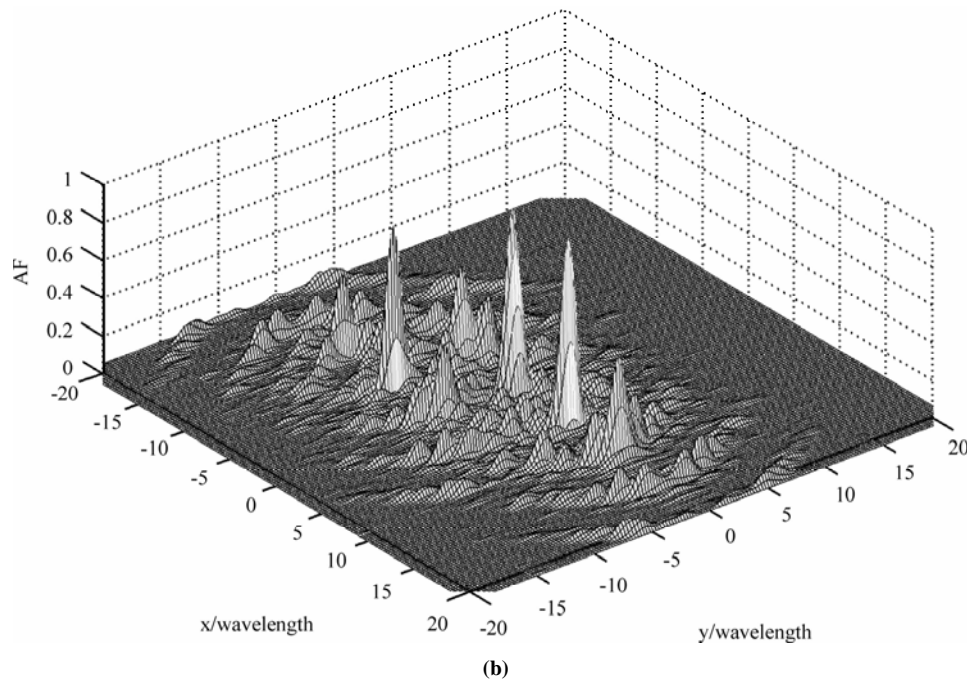
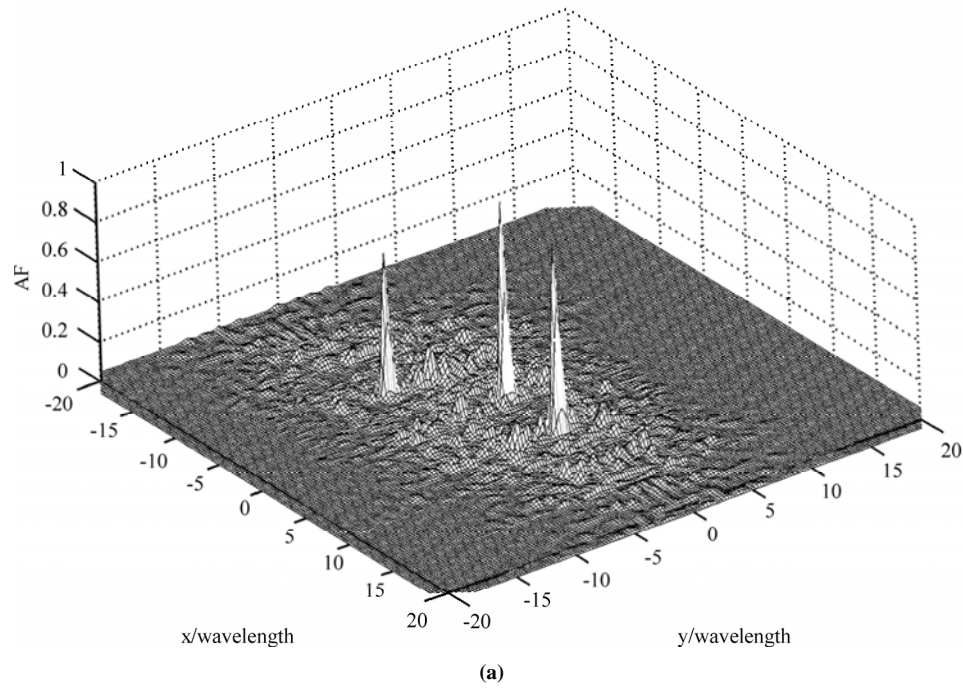
The mesh plots of the AF are showed in Figure 2 and more details on resolutions and sidelobe characteristics are given in Figure 3 and Table 1. As is obvious from both Figure 3 and Table 1, the resolutions of MIMO-HFSWR are at the level of one wavelength for both



**Figure 3. Resolution and sidelobe performances of sparse frequency waveforms (solid line) and single LFM signal (dotted line) along (a) x-axis and (b) y-axis.**

**Table 1. Resolution and sidelobe characters of different waveforms.**

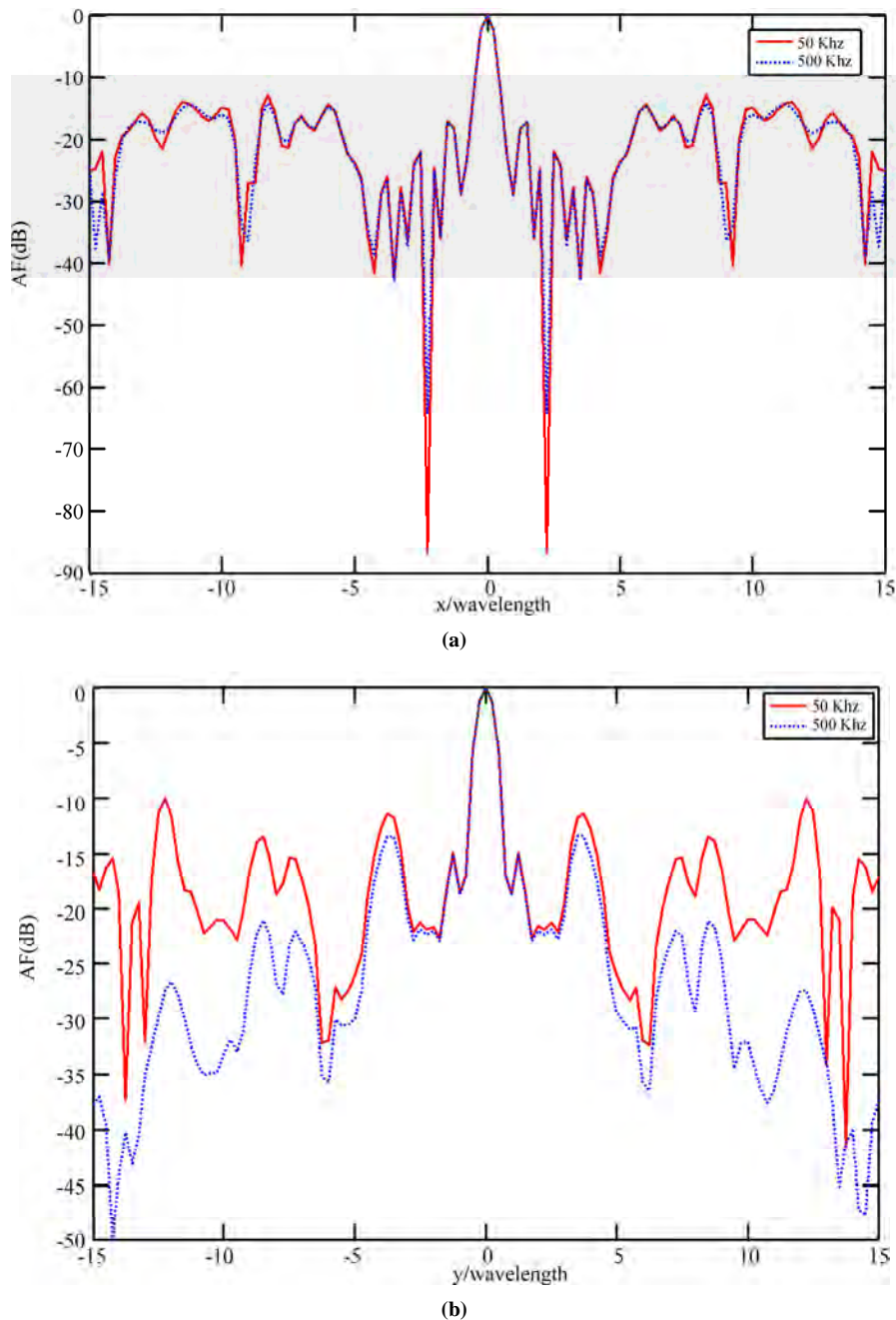
Item	Sparse frequency waveforms	LFM
Resolution of x	0.31 wavelength	0.58 wavelength
PSL-x	-14.9 dB	-3.2 dB
Resolution of y	0.9 wavelength	1.6 wavelength
PSL-y	-13.5 dB	-16.5 dB

**Figure 4. Resolution capacity in 3 point targets case (a) Sparse frequency waveforms (b) Single LFM waveform.**

sparse frequency waveforms and single LFM waveform. And the resolutions of sparse frequency waveforms are even better than those of common LFM signals. This is a great improvement for azimuth resolution and even for range resolution from conventional several kilometers to several tens meters. Meanwhile, as HFSWR always works in a highly congested spectrum environment, the sparse frequency waveform approach can provide better flexibility on choosing available channels than wave-

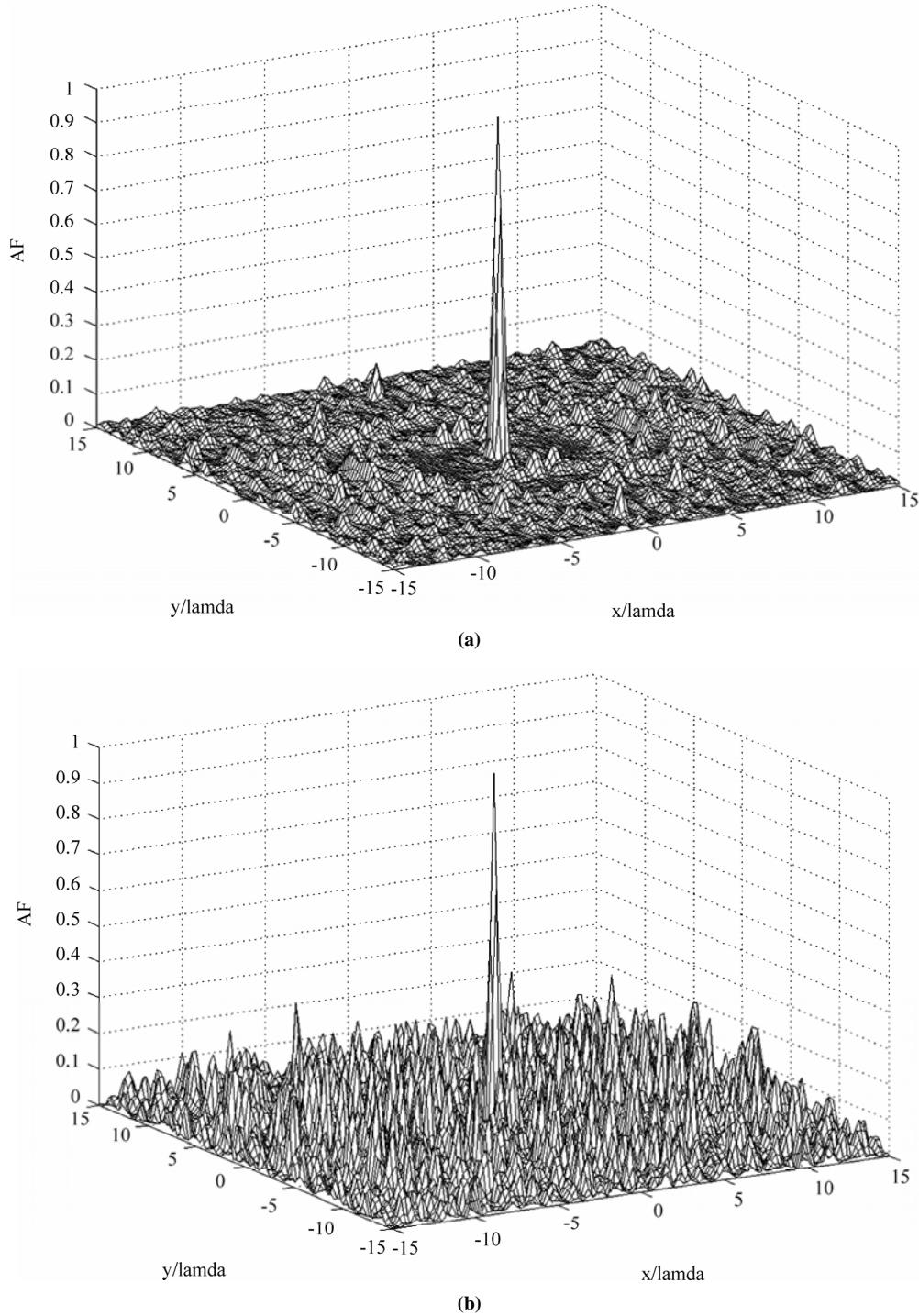
forms confined in only one channel. The sidelobes of x-axis and y-axis are well below  $-13$  and  $-14$  dB for sparse frequency waveforms, respectively, which is a significant improvement compared with the side lobe level from the single LFM waveform within the same channel. It demonstrates that the sidelobe levels can be suppressed by frequency diversity in random arrays [13]. This is another advantage of sparse frequency waveform.

Multiple targets case are illustrated in Figure 4, where



**Figure 5. Bandwidth effect on resolution and sidelobe performance along (a) x-axis and (b) y-axis. Solid line is associated with bandwidth 50 KHz, while the dotted line is 500KHz.**

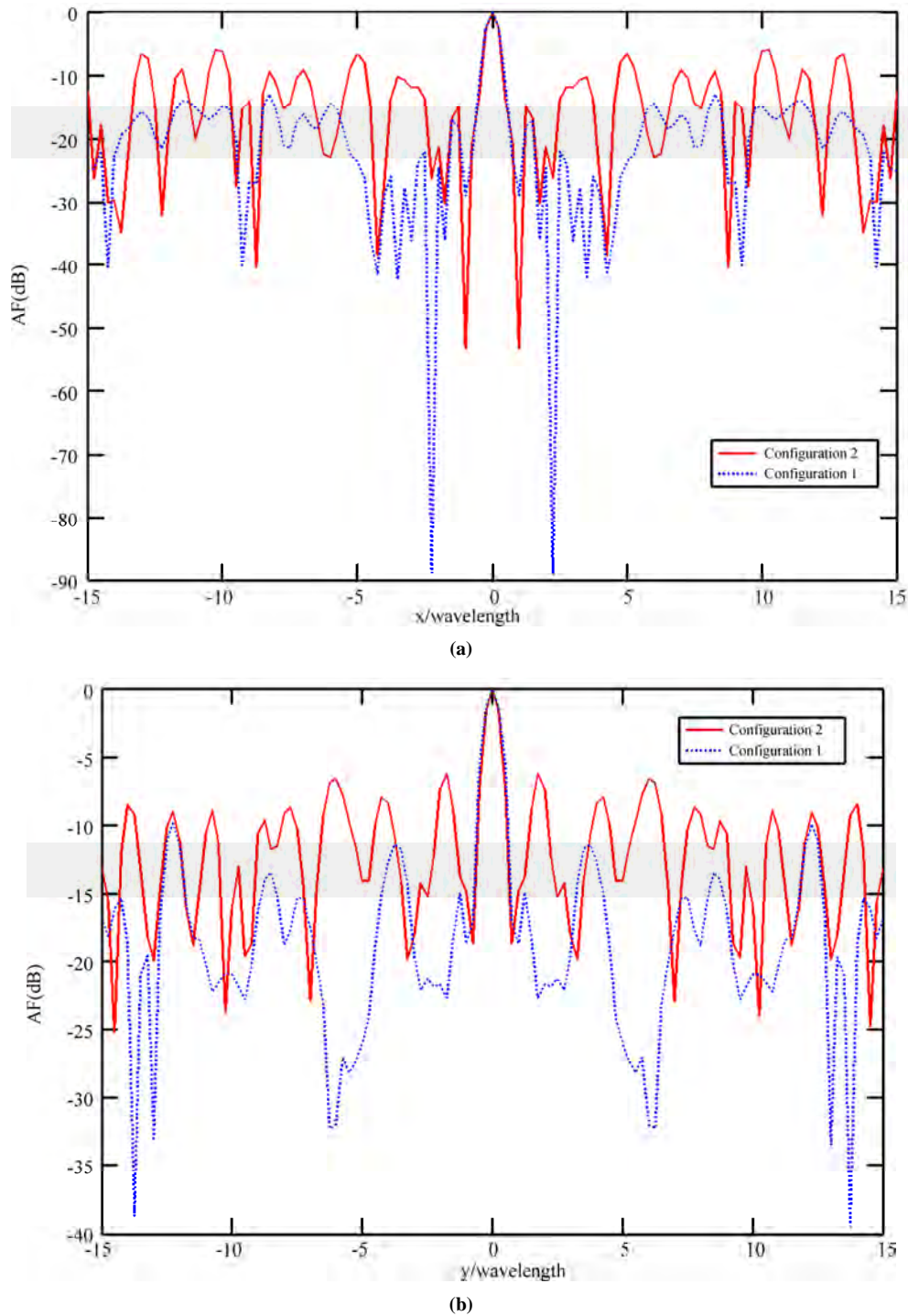




**Figure 6.** Ambiguity functions of (a) Configuration 1 and (b) Configuration 2 in case 1 with velocity  $[v_x, v_y] = [500, 500]$  m/s. Configuration 1 (9×9), Configuration 2 (5×5), both evenly distributed in  $(-\pi/4, \pi/4)$ .

three targets are located in  $[0, 0]$ ,  $[0, 10]$ , and  $[-10, -10]$ . The coordinate system is expressed in multiples of wavelength. As we can see, by using sparse frequency waveform set, the system can better distinguish different targets than by using waveform set in the same channel.

Bandwidth effect is illustrated in Figure 5. We take a set of sparse waveforms like that mentioned above. The difference is that the bandwidth is 50 KHz for each waveform. From Figure 5 we can see that even with smaller bandwidth, the resolution capacity is not much impacted.



**Figure 7. Resolution and sidelobe performance along (a)  $x$ -axis and (b)  $y$ -axis in case 1 with velocity  $[v_x, v_y] = [500, 500]$  m/s. Configuration 1 (9×9), Configuration 2 (5×5), both evenly distributed in  $(-\pi/4, \pi/4)$ .**

However, the PSL performance is deteriorated. The PSL in  $x$ -axis is about  $-12.9$  dB and is about  $-10$  dB in  $y$ -axis for this waveform set. Thus, we can see that the larger the bandwidth adopted, the lower the sidelobes in both  $x$ -axis and  $y$ -axis.

In this case study, the simulation results demonstrate that MIMO-HFSWR with sparse frequency waveform has superior resolution than conventional HFSWR in both range and downrange domain. Sidelobe levels can be suppressed by using sparse frequency waveforms.

Further work will be focused on the suppression of sidelobe levels by waveforms with effective frequency diversity scheme.

#### 4. Doppler and Geometry Factor

As HFSWR is always operated in Doppler circumstances, the AF with Doppler effects should be further investigated. Meanwhile, unlike monostatic radar the distributed MIMO radar is confined by the geometry configuration. Thus the geometry factor should also be investigated. Two cases are given below to investigate the Doppler and configuration effect.

In Case 1, we study the configuration effect. Target of this case is with velocity of  $[v_x, v_y] = [500, 500]$  m/s. This is a high velocity target case corresponding to air targets. There are two configurations. For Configuration 1, in the region of  $(-\pi/4, \pi/4)$  there are nine transmitters and nine receivers, evenly distributed. The transmitting waveforms are defined as those in Section 3 except that each one has 10 kHz bandwidth. For practical HFSWR application, only a limited number of continuous clear channels with bandwidth of a few kilo-Hertz in the 3-30 MHz high frequency band can be found and used at a time when interference is considered [1,2]. Thus, 10 kHz bandwidth is used for a much more similitude in real condition of the HFSWR system. For Configuration 2, in

the same region of  $(-\pi/4, \pi/4)$  there are five transmitters and five receivers, evenly distributed. The waveforms are the first five used in Configuration 1 of Case 1. Thus, the total spectra employed by these two configurations are the same. As illustrated in Figure 6, Figure 7, both configurations in this case show high resolution capabilities. However, Configuration 1 with more transmit-receive pairs shows better sidelobe performance in both x-axis and y-axis. The PSL in x-axis is about -12 dB and is about -10.5 dB in y-axis for Configuration 1. Based on our numerous simulation experiments, it is found that as more pairs of transmitter and receiver are set in a much wider spatial region, the resolutions can be slightly improved and the PSLs of y-axis and x-axis can be further reduced. However, systematic study on the PSLs reduction through geometry optimization will be explored in the future.

In Case 2: we have four velocity settings like  $[0, 0]$  m/s,  $[100, -100]$  m/s,  $[-10, 5]$  m/s, and  $[500, 500]$  m/s. The geometry configuration in Case 2 is the same as Configuration 1 in Case 1. We also take the waveform set of Configuration 1 of Case 1 in this case study. Figure 8 shows the results of Case 2. We can see from Figure 8 that the proposed system shows similar characteristics in different Doppler context, which means the resolution and PSL performance are both insensitive to Doppler frequency. Thus, for both high speed air targets and low

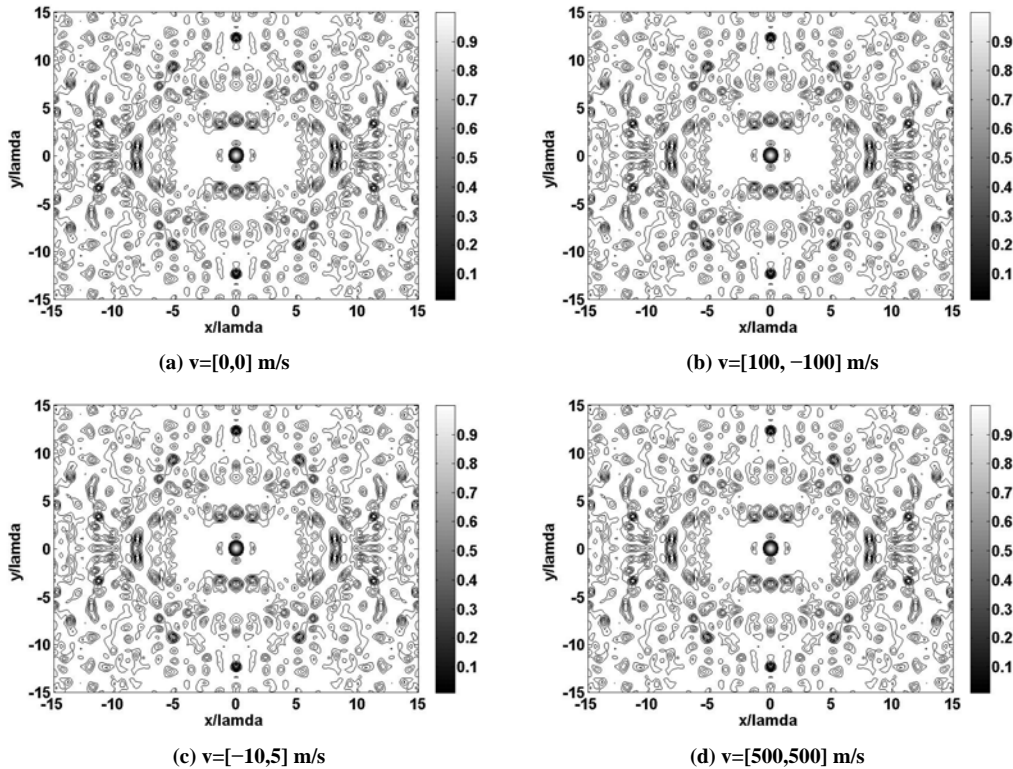


Figure 8. Ambiguity functions of different velocity with Configuration 1.

velocity surface targets, the proposed system can also have high resolution performance.

## 5. Conclusions

In this paper, the concept of distributed MIMO-HFSWR radar transmitting sparse frequency waveforms is proposed. The AF of this proposed system is derived in detail. Potential advantages of the proposed system on resolution capacity and PSL performance are assessed through AF analysis and simulations. The impacts of Doppler effects and the geometry configuration factor are also studied. It has been found that the system has several distinguished characteristics. Firstly, the range resolution and the azimuth resolution can be improved to the level of one wavelength, namely, only tens meters and the PSL is reduced to a much lower level with sparse frequency waveforms. Meanwhile, the resolutions are not restricted by individual bandwidth while the PSL can benefit from large bandwidth. Secondly, the performance of fine resolution and low PSL are insensitive to the Doppler effects. Thus, for both high speed air and low velocity surface targets, the proposed system also has high performance. Thirdly, the resolution capacity and PSL performance can be optimized through geometry configuration optimization. In addition, multistatic configuration provides large flexibility to find a proper place to locate the radar transmitters and receivers; by using sparse frequency waveforms, it is much easier to find more available channels in different locations thus the co-channel interference can be avoided and the performance can be further improved. Further studies will be conducted on the surveillance strategy and high quality waveforms with better AF performance. Meanwhile, synchronization problem should also be paid special attention to so that coherent processing can be conducted perfectly.

## 6. References

- [1] H. W. H. Leong and B. Dawe, "Channel availability for east coast high frequency surface wave radar systems," Defence R&D Canada, Technical Report, DREO TR 2001-104, November 2001.
- [2] R. J. Riddolls, "A Canadian perspective on high frequency over-the-horizon radar," Defence R&D Canada, Technical Report, DREO TR 2006-285, December 2006.
- [3] D. W. Bliss and K. W. Forsythe, "Multiple-input multiple-output (MIMO) radar and imaging: Degrees of freedom and resolution," Proceedings of the 37th Asilomar Conference on Signal, Systems and Computers, pp. 54–59, November 2003.
- [4] J. Li, "MIMO radar: Diversity means superiority," Proceedings of the 14th Annual Adaptive Sensor Array Processing Workshop-2006, June 6–7, 2006.
- [5] S. Peter, J. Li, and Y. Xie, "On probing single design for MIMO radar," IEEE Transactions on Signal Processing, Vol. 55, No. 8, pp. 4151–4161, August 2007.
- [6] G. San Antonio, D. R. Fuhrmann, and F. C. Robey, "MIMO radar ambiguity function," in IEEE Journal of Selected Topics in Signal Processing, Vol. 1, No. 1, pp. 167–177, June 2007.
- [7] E. Fisher, A. Haimovich, R. Blum, L. Cimini, D. Chizhik, and R. Valenzuela, "Spatial diversity in radars—models and detection performance," IEEE Transactions on Signal Processing, Vol. 20, No. 3, pp. 823–838, March 2006.
- [8] N. H. Lehmann, A. M. Haimovich, R. S. Blum, and L. Cimini, "High resolution capabilities of MIMO radar," Proceedings of the 40th Asilomar Conference on Signal, Systems and Computers, pp. 25–30, November 2006.
- [9] D. R. Kirk, J. S. Bergin, P. M. Techau, and J. E. Don Carlos, "Multi-static coherent sparse aperture approach to precision target detection and estimation," Proceedings of 2005 IEEE Radar Conference, pp. 579–584, May 2005.
- [10] G. H. Wang, W. X. Liu, and Y. L. Lu, "Sparse frequency transmit waveform design with soft power constraint by using PSO algorithm," Proceedings of 2008 IEEE Radar Conference, pp. 1–6, May 2008.
- [11] W. X. Liu, Y. L. Lu, and M. Leisturgie, "Optimal sparse waveform design for HFSWR system," Proceedings of 2007 International Waveform Diversity and Design Conference, pp. 127–130, May 2007.
- [12] N. Levanon and E. Mosezen, Radar Signals, John Wiley & Sons, New Jersey, 2004.
- [13] B. D. Steinberg and E. H. Attia, "Sidelobe reduction of random arrays by element position and frequency diversity," IEEE Transactions on Signal Processing, Vol. 31, No. 6, pp. 922–930, November 1983.

# ContSteg: Contourlet-Based Steganography Method

Hedieh SAJEDI, Mansour JAMZAD

Computer Engineering Department, Sharif University of Technology, Tehran, Iran

E-mail: A\_sajedi@ce.sharif.edu, Jamzad@sharif.edu

Received April 26, 2009; revised May 20, 2009; accepted May 25, 2009

## Abstract

A category of techniques for secret data communication called steganography hides data in multimedia mediums. It involves embedding secret data into a cover-medium by means of small perceptible and statistical degradation. In this paper, a new adaptive steganography method based on contourlet transform is presented that provides large embedding capacity. We called the proposed method ContSteg. In contourlet decomposition of an image, edges are represented by the coefficients with large magnitudes. In ContSteg, these coefficients are considered for data embedding because human eyes are less sensitive in edgy and non-smooth regions of images. For embedding the secret data, contourlet subbands are divided into  $4 \times 4$  blocks. Each bit of secret data is hidden by exchanging the value of two coefficients in a block of contourlet coefficients. According to the experimental results, the proposed method is capable of providing a larger embedding capacity without causing noticeable distortions of stego-images in comparison with a similar wavelet-based steganography approach. The result of examining the proposed method with two of the most powerful steganalysis algorithms show that we could successfully embed data in cover-images with the average embedding capacity of 0.05 bits per pixel.

**Keywords:** Information Hiding, Steganography, Steganalysis, Contourlet Transform

## 1. Introduction

Steganography methods hide the secret data in a cover carrier so that the existence of the embedded data is undetectable. The cover carrier can be different kinds of digital media such as text, image, audio, and video [1]. In a successful steganography method the carrier medium does not attract attentions. The security of the steganography methods is mostly influenced by the kind of cover media, the method for selection of places within the cover that might be modified, the type of embedding operation, and the number of embedding changes that is a quantity related to the length of the embedded data.

The aim of the steganography methods is to communicate securely in a completely undetectable manner. As the steganography techniques progress, there is an increased interest in steganalysis algorithms which their main goal is detecting the presence of hidden data.

Many steganography methods have been proposed and several stego-products have been developed (e.g., EzStego [2]) in which an innocuous-looking image is used as the cover-image to conceal the secret data. In these methods, the secret data is embedded into the cover-image

by modifying the cover-image to form a stego-image.

Some image hiding systems use uncompressed images (e.g., BMP) or lossless compressed images (e.g., GIF) as cover-images. These images potentially contain visual redundancy so that they can provide large capacity to hide secret data. For reducing transmission bandwidth and storing space, the JPEG is currently the most common format for images that are used on the Internet. Therefore, embedding techniques in Discrete Cosine Transform (DCT) domain are popular because of the large usage of JPEG images. Although modifications of properly selected DCT coefficients during embedding process will not cause noticeable visual artifacts, nevertheless they cause detectable statistical degradations. Various steganography methods like F5 [3], Outguess [4], Model-based (MB) [5], Perturbed Quantization (PQ) [6], and YASS [7] have been proposed with the purpose of minimizing the statistical artifacts which are produced by modifications of DCT coefficients.

On the other hand, some steganography methods based on wavelet transform have been presented. In [8], a steganography method based on wavelet and modulus function is proposed. In this method, the capacity of a



cover-image is determined considering the number of wavelet coefficients with larger magnitude.

Embedding data in adaptively selected parts of cover-images such as regions having edges and texture enhances the security of stego-images [9]. An adaptive steganography method attempts to provide secure embedding by ensuring that the changes introduced into the cover-images remain consistent with natural properties of them. Since human eyes are less sensitive in edgy and non-smooth regions of images, modifications in these parts of cover-images are less detectable.

In [10] we proposed a new steganography method that embeds secret data in contourlet coefficients of images. In this paper, we describe the method introduced in [10] with more details and complete our experiments with a larger image database. In this paper, we introduce ContSteg, which is a method based on contourlet transform for hiding data in images. In ContSteg, contourlet transform is applied to capture significant image coefficients across spatial and directional resolutions. Multiresolution flexibility, local and directional image expansion in the contourlet image representation, allow for easy subband processing [11]. To increase the embedding capacity and quality of stego-images compared to previous methods, we embed the secret data in proper contourlet coefficients of the cover-image. The embedding algorithm takes advantage of adaptive methods by embedding data in non-smooth regions of cover images. In this way, the visual degradation caused by the steganography method can be mitigated because the secret data is embedded in higher contourlet coefficients in edgy and non-smooth areas that can visually hide this information better [12]. The embedding process is carried on by changing the

value of two contourlet coefficients to hide one bit of secret data.

The experimental results illustrated that the proposed method can hide much more data while maintaining a good visual quality of stego-images compared to the similar wavelet-based steganography methods. We verified that by employing two well-known and efficient steganalysis methods. They could not discriminate between clean and stego-images reliably.

The rest of this paper is organized as follows. In Section 2, we introduce the proposed steganography method, ContSteg, and discuss the main characteristics of contourlet transform. Performance of the presented method is analyzed in Section 3 and finally, we conclude this paper in Section 4.

## 2. ContSteg

Using suitable representation domain and proper coefficients to embed data, can result in stego-images with higher quality. Consequently, higher embedding capacity and enhanced security are provided. Accordingly, in this paper, a new method is proposed which is called ContSteg. It takes advantage of a multiscale framework and its directionality to extract the appropriate places of an image to hide data. ContSteg like other steganography methods consists of an embedding process and an extraction process. Figure 1 shows the block diagram of embedding and extraction processes of ContSteg. The details of these processes are described in the following subsections.

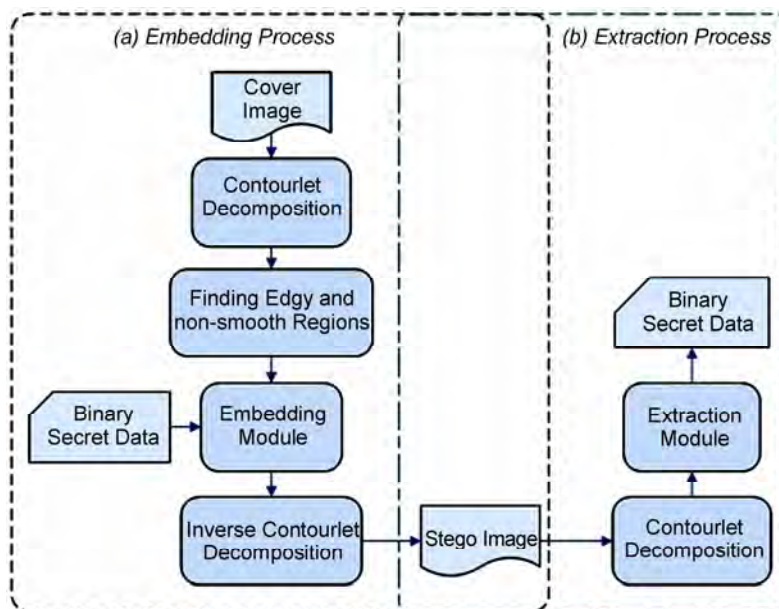
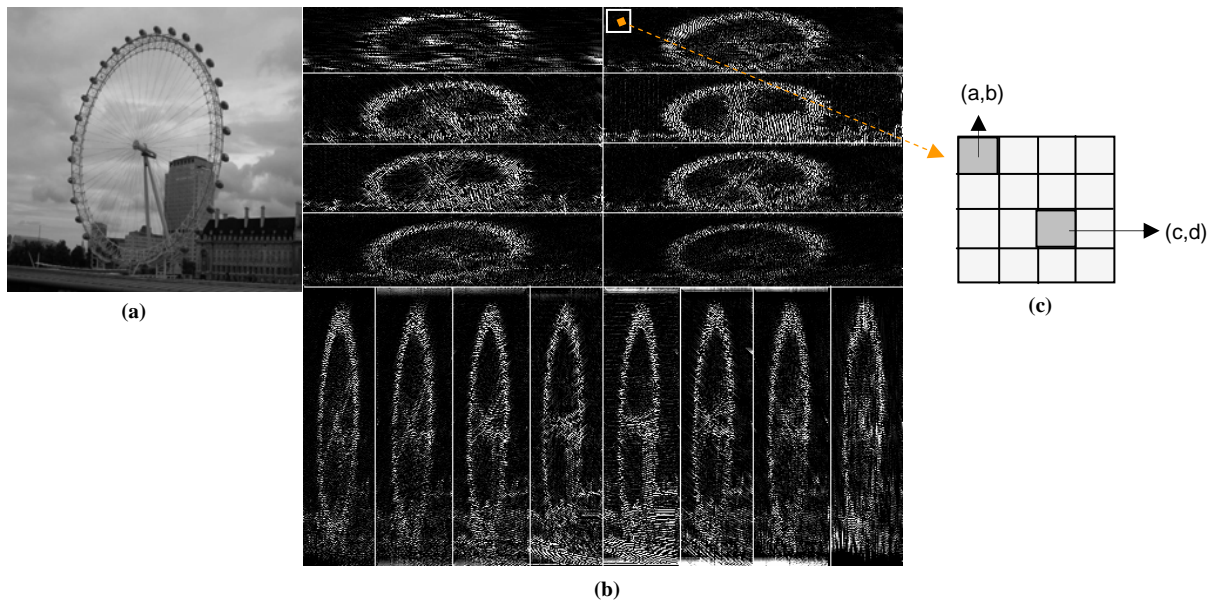


Figure 1. The block diagram of ContSteg steganography method, (a) Embedding process, (b) Extraction process.



**Figure 2.** Embedding data in contourlet coefficients of an image, (a) Original input image, (b) Visualization of contourlet decomposition of an image into one pyramidal level and sixteen directional subbands, (c) A  $4 \times 4$  block of contourlet coefficients and the place of two coefficients for embedding.

## 2.1. Hiding Data in Contourlet Coefficients

Contourlet transform is one of several transforms developed in recent years, aimed at improving the representation sparsity of images over the wavelet transform. The main feature of this transform is the potential to handle 2-D singularities efficiently, i.e. edges, unlike wavelet, which can deal with point (i.e. 1-D) singularities exclusively [13]. Contourlet transform is a directional extension of wavelet transform that fixes the wavelet subband-mixing problem and improves its directionality. Two-dimensional wavelet transform produces one approximation subband, and three details subbands, corresponding to the horizontal, vertical, and diagonal directions. The diagonal subband mixes the directional information oriented at  $45^\circ$  and  $135^\circ$ . The main idea of contourlet is to find some directional extensions to divide further each detail subband of the wavelet into a number of directions. This transform is based on a double filter bank structure by combining the Laplacian pyramid with a directional filter bank [14]. Figure 2 shows an image that is decomposed into one pyramidal level and sixteen directional subbands (higher coefficients are colored white).

Because of the subband-mixing problem in wavelet transform, manipulating one coefficient in diagonal subband affects the value of other relevant coefficients in other directions. We used the effectiveness of contourlet transform in image decomposition to separate directions. Hence, manipulating the value of a coefficient in the contourlet subbands has less effect in the quality of the

image than changing a coefficient in wavelet subbands. Furthermore, most of the current existing steganalysis algorithms are limited to the domain of spatial, wavelet, and DCT transform. Therefore, distinguishing cover-images from stego-images (constructed by embedding data into their contourlet coefficients) is not easy by these steganalysis algorithms. Accordingly, considering the fact that higher embedding efficiency translates into better steganographic security, more secure stego-images are achieved using the proposed method.

## 2.2. Embedding Process

The embedding process is done in the following steps:

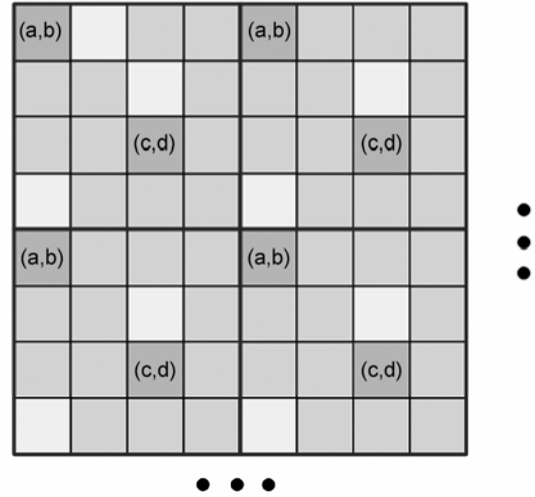
- Step 1:* The cover-image is decomposed with one pyramidal level and sixteen directional contourlet transform.
- Step 2:* The regions of the subbands in which the data can be embedded are identified. Then the embedding process determines higher contourlet coefficients in these regions that can be used for embedding.
- Step 3:* According to Kerckhoffs' principle [15], the embedding algorithm is supposed to be known to the public. Therefore, the embedding process may use an embedding key so that only the legal user can successfully extract the embedded data by using the corresponding extraction key in the extraction process. Accordingly, a key that is a seed for generating a random sequence is considered to provide the embedding location addresses of  $4 \times 4$  blocks.

**Step 4:** In this step, the embedding module is activated. The place of two coefficients in each block are chosen by the embedding module and agreed upon by both send and receive parties. These two coefficients are suitable for embedding if both of them belong to the higher coefficients set. The embedding module hides each bit of the secret data by comparing and if needed exchanging the values of two contourlet coefficients in non-smooth regions of the image. We use two coefficients that are shown in Figure 2(c). A  $4 \times 4$  block encodes bit 1 if its  $coefficient(a,b) \geq coefficient(c,d)$  and bit 0 otherwise. Two coefficients are swapped if their values do not match with the bit to be encoded. Since the JPEG compression, rounding in computation, and non-orthogonality of contourlet transform can affect the relative size of the coefficients, the embedding module ensures that  $|coefficient(a,b) - coefficient(c,d)| > t$ , where  $t$  is a value that represents the tradeoff between image quality and hidden data retrieval error rate. We set  $t = 2$  experimentally. Due to the cases we mentioned before, manipulating the value of coefficients may cause loss of the embedded data in inverse contourlet transform. In addition, it may affect the value of neighborhood coefficients and thus the embedded data in such neighborhood may be lost. To maintain a high level of similarity between the original clean and stego-images, and to have minimum loss in extracted data, each candidate coefficient for embedding should have a distance from other candidate coefficients. Considering these properties, we embed each bit in coefficient block of size  $4 \times 4$ . In this fashion, a candidate coefficient has the least closeness to other candidates. Figure 3 shows a part of a contourlet subband, which has some  $4 \times 4$  blocks. As the figure shows, candidate coefficients for embedding are considered far from other candidates.

### 2.3. Extraction Process

The stego-key used in the embedding process should be shared by both the sender and receiver so that the embedded data can be extracted by a legal receiver. The extraction module consists of the following steps:

- Step 1:** Decompose stego-image with a one level contourlet transform.
- Step 2:** Recognize higher contourlet coefficients.
- Step 3:** Form the random sequence by using the same key as the sender has used.
- Step 4:** Retrieve the embedded data by comparing  $coefficient(a,b)$  and  $coefficient(c,d)$  in each  $4 \times 4$  coefficient block. If  $coefficient(a,b) \geq coefficient$



**Figure 3.** A part of a contourlet subband with some  $4 \times 4$  blocks. Candidate coefficients (shown in dark gray) for embedding are at least one pixel apart from other candidates.

$(c,d)$ , the hidden bit is 1 and it is 0 otherwise.

Figure 1(b) shows the block diagram of the extraction process of ContSteg.

## 3. Experiments

We did different experiments to assess the efficiency of the proposed method. We collected 1000 images from some typical images and some random ones from Washington University image database [16]. All images were converted to grayscale and cropped to size of  $512 \times 512$ . The JPEG quality factor of images is 75. To obtain a stego-dataset, for each cover-image a random binary data was embedded using ContSteg. Therefore, in our database we have 2000 images, 1000 cover-images, and 1000 stego-images.

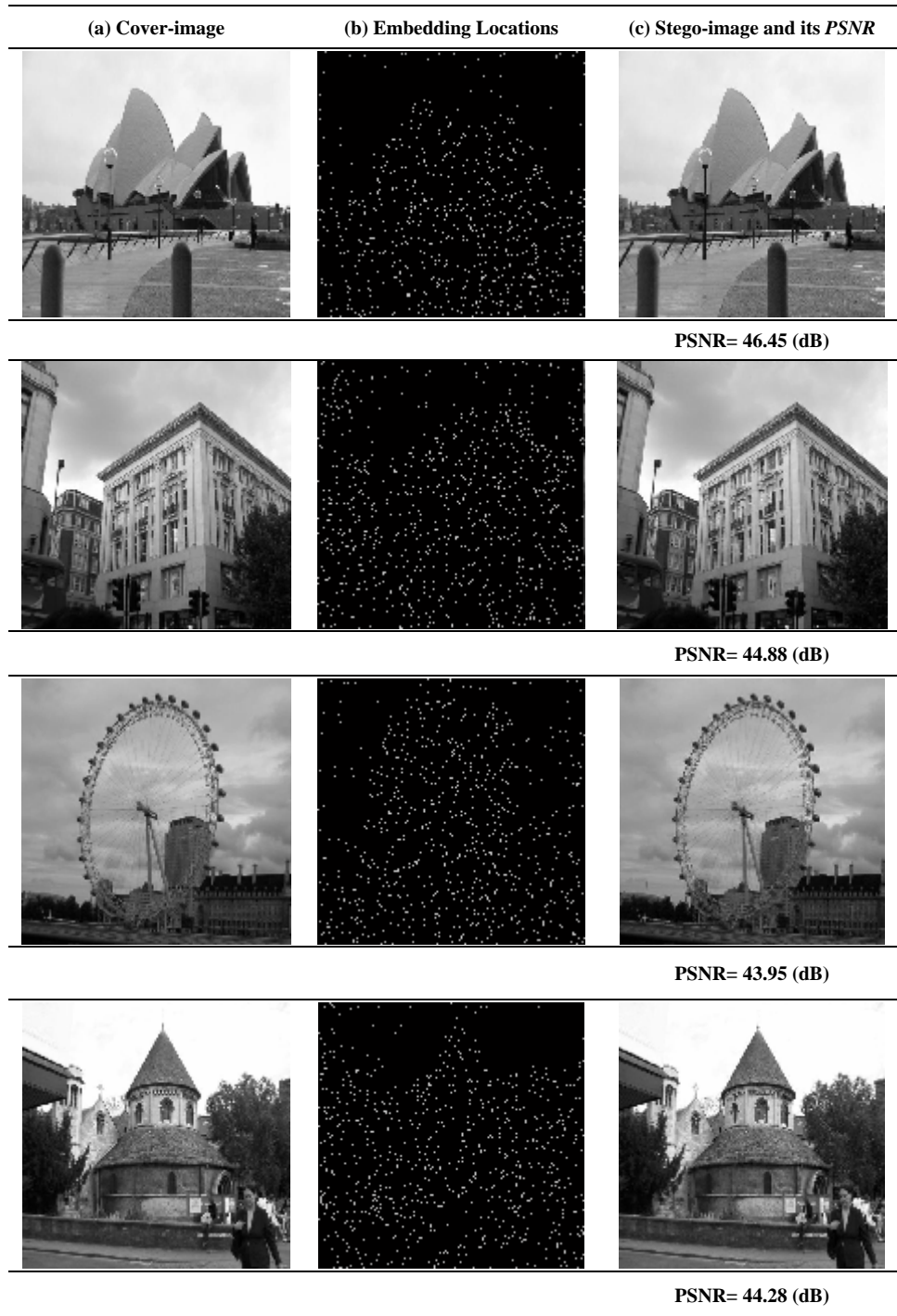
### 3.1. Efficiency of ContSteg

In this experiment, we assess the efficiency of ContSteg in terms of quality of stego-images and embedding rate of ContSteg.

#### 3.1.1. Calculation of Embedding Rate

In the proposed method, the desired frequency partitioning for a  $N \times N$  size image by contourlet transform contains of sixteen directional subbands of size  $N/8 \times N/2$  in first level of decomposition. By embedding one bit of secret data in each  $4 \times 4$  block of all subbands, the embedding capacity of an image will be  $(N \times N)/16$ . If  $C$  percents of coefficients are used for embedding, then the embedding rate is  $C/16$  bits per pixel. In most of the steganography methods based on wavelet transform, approximation subband is not used for embedding. Because





**Figure 4.** Computing the quality of stego-images, (a) Cover-image, (b) Proper locations for embedding are colored white, (c) Stego-image with its PSNR.

changing the coefficients in approximation subband imposes a large distortion in the stego images. Hence, in this case the embedding rate should be very low. Therefore, for a  $N \times N$  image, in the first level of decomposition, the number of contourlet coefficients is  $(N \times N)/4$  more than wavelet coefficients. Therefore, more embedding

rate can be archived in this domain. For a  $512 \times 512$  size image, the number of contourlet coefficients is 262144. This number is equal to the number of image pixels. If we keep 50 percent of higher coefficients, we have 131072 coefficients. If one bit is embedded in each  $4 \times 4$  block, the maximum rate for embedding is about 0.03

bits per pixel. By this configuration, in the best condition (since we embed only in coefficients with higher amplitude, a block is proper for embedding if it has coefficients with higher amplitude) we can embed 8192 bits of data in the mentioned image. Using greater percent ( $C > 50$ ) of coefficients with higher amplitude for embedding provides higher embedding rate ( $> 8192$ ).

### 3.1.2. Computing the Quality of Stego-Images

In this evaluation, we consider perceived quality of stego-images. Figure 4 shows some cover-images, the locations to embed data and the stego-images after embedding 5600 bits. The results show that the quality of stego-images is high, and unintended observers cannot be aware of the existence of hidden data in it. The imperceptibility is evaluated by the objective quality measurement *PSNR* (peak signal to noise ratio) [17]:

$$PSNR = 10 \times \log \left( \frac{255^2}{MSE} \right) \quad (1)$$

where *MSE* represents the mean square error between the cover-image  $x$  and the stego-image  $y$  both of size  $512 \times 512$ .

$$MSE = \left( \frac{1}{512 \times 512} \right) \sum_{i=1}^{512} \sum_{j=1}^{512} (x_{ij} - y_{ij})^2 \quad (2)$$

Figure 5 shows the average *PSNR* for images of size  $512 \times 512$  after embedding the secret data of size 3000 to 12000 bits in wavelet and contourlet coefficients of images. In this figure, the points on the curve correspond to the average *PSNR* of stego-images in the database with certain payloads. For example, for payload of 3000 bits some stego images in our database has *PSNR* above 45 (dB) and some other have *PSNR* below 35 (dB) but averagely *PSNR* is about 38.8 (dB). The embedding and extraction processes in wavelet and contourlet domains are the same. The results show that embedding in con-

tourlet transform domain increases the quality of stego-images.

### 3.2. Protection against JPEG Compression

Due to the rounding in computation, and non-orthogonality of wavelet and contourlet, embedding methods in both of these domains have less than 1% loss of the secret data in the worst case. For lossless data recovery, we have to use a redundancy factor in an error correction framework. Table 1 shows the evaluation of proposed steganography technique against JPEG compression. As we see, the proposed method has not a good robustness against compression but with the cost of lower quality stego-images (e.g. using hamming code algorithm that makes the secret data secure with added redundancy), higher robustness against compression can be achieved.

### 3.3. Steganalysis Results

Wavelet-based steganalysis (WBS) [18], and Feature-based steganalysis (FBS) [19], and Contourlet-based (CBS) [20] methods are used to evaluate the security of ContSteg. In WBS, a Fisher Linear Discriminator (FLD) and in FBS and CBS, a nonlinear Support Vector Machine (SVM) is trained to discriminate between clean and stego-images. 1200 images (600 cover and 600 stego images) from database were chosen randomly for testing, while the remaining 800 images were used for training. This partitioning was repeated ten times, with different random subsets used for training and testing each time. The average of detection accuracy is shown in Table 2. The accuracy is the average of true detection of both stego and clean-images. As can be seen, the detection accuracy is about 50% and the proposed method with payload of approximately 0.05 bits per pixel cannot be reliably detected by the applied steganalyzers.

It is shown in [21] that the average embedding capacity of existing steganography methods for grayscale JPEG

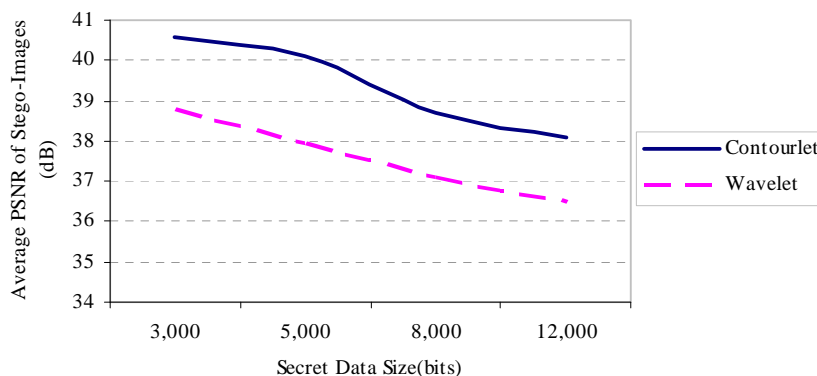


Figure 5. Comparing the quality of stego-images produced by wavelet-based and contourlet-based steganography methods.

**Table 1. Retrieval error rate of hidden data after JPEG compression.**

Secret Data Size (bits)	Quality Factor	Retrieval Error Rate (%)
5,000	90 , 70 , 50	10 , 14 , 20
10,000	90 , 70 , 50	13 , 18 , 26

**Table 2. Accuracy of WBS, FBS, and CBS steganalysis methods on detection of stego-images produced by ContSteg.**

Secret Data Size (bits)	Steganalysis Method	Average Detection Accuracy (%)
5,000	WBS	51
	FBS	53
	CBS	59
10,000	WBS	53
	FBS	54
	CBS	63
15,000	WBS	58
	FBS	61
	CBS	68

images with quality factor of 70 is approximately 0.05 bits per non-zero AC DCT coefficient. For a  $512 \times 512$  image, 4096 blocks of size  $8 \times 8$  is existed. Usually 20 AC DCT coefficients are considered non-zero. Therefore, we have  $4096 \times 20 = 81920$  non-zero coefficients. Hence, the capacity is  $81920 \times 0.05 = 4096$  which is  $4096 / (512 \times 512) = 0.015$  bits per pixel. We see that our proposed method has higher embedding capacity.

#### 4. Conclusions

Steganography that is a branch of information hiding technology aims to hide a secret data securely in a cover media for transmission. Embedding rate and stego-image quality are two important criteria in evaluating a steganography method. In this paper, a new secure and adaptive steganography is presented which is called ContSteg. It embeds a secret data in contourlet transform coefficients of an image. Since embedding data in non-smooth and edgy regions of the image causes less detectability, these regions of the image are identified in contourlet domain and the secret data is embedded in the corresponding coefficients. According to the experimental results, in comparison with wavelet domain approach, the proposed steganography method increases embedding rate and image quality of the stego-images by hiding the secret data in contourlet coefficients corresponding to high frequencies. The results of our experiments show that employing two of powerful steganalyzers on stego-images produced by our method, they could not

discriminate between stego and clean-images reliably. In general, ContSteg is a secure steganography method that provides high embedding capacity and high image quality.

#### 5. References

- [1] C. Liu and S. Liao, "High-performance JPEG steganography using complementary embedding strategy," *Pattern Recognition*, Vol. 41, pp. 2945–2955, 2008.
- [2] EzStego, <http://www.securityfocus.com/tools/586>.
- [3] A. Westfeld, "F5-a steganographic algorithm: High capacity despite better steganalysis," *Proceeding of 4th International Workshop on Information Hiding*, 2001.
- [4] N. Provos, "Defending against statistical steganalysis," *Proceeding of 10th USENIX Security Symposium*, pp. 323–336, 2001.
- [5] P. Sallee, "Model-based steganography," *Proceeding of International Workshop on Digital Watermarking*, Seoul, Korea, 2003.
- [6] J. Fridrich, M. Goljan, and D. Soukal, "Perturbed quantization steganography with wet paper codes," *Proceeding of ACM Multimedia Workshop*, Germany, 2004.
- [7] K. Solanki, A. Sarkar, and B. S. Manjunath, "YASS: Yet another steganographic scheme that resists blind steganalysis," *Proceeding of 9th International Workshop on Information Hiding*, June 2007.
- [8] K. Zhiwei, L. Jing, and H. Yigang, "Steganography based on wavelet transform and modulus function," *Journal of*

- Systems Engineering and Electronics, Vol. 18, No. 3, pp. 628–632, 2007.
- [9] J. Fridrich and R. Du, “Secure steganographic methods for palette images,” *Proceeding of 2nd International Information Hiding Workshop. LNCS*, Vol. 1768, pp. 47–60, 2000.
  - [10] H. Sajedi and M. Jamzad, “Adaptive steganography method based on contourlet transform,” *Proceedings of 9th International Conference on Signal Processing (ICSP’08)*, October 26–29, 2008.
  - [11] M. Do and M. Vetterli, “Contourlets: A directional multiresolution image representation,” *Proceedings of ICIP*, 2002.
  - [12] N. Kaewkamnerd and K. R. Rao, “Wavelet based image adaptive watermarking scheme,” *Electronic Letters*, Vol. 36, pp. 312–313, 2000.
  - [13] B. Matalon, M. Elad, and M. Zibulevsky, “Image denoising with the contourlet transform,” *Proceeding of SPIE Conference Wavelets*, 2005.
  - [14] Y. Lu and M. N. Do, “A directional extension for multi-dimensional wavelet transforms,” *IP EDICS: 2-WAVP (Wavelets and Multiresolution Processing)*, 2005.
  - [15] J. Seberry and J. Pieprzyk, “CRYPTOGRAPHY: An introduction to computer security,” Prentice-Hall, New York, 1989.
  - [16] <http://www.cs.washington.edu/research/imagetdatabase>.
  - [17] A. K. Jain, “Fundamentals of digital image processing,” Prentice-Hall, New Jersey, 1989.
  - [18] S. Lyu and H. Farid, “Detecting hidden messages using higher-order statistics and support vector machines,” *Proceeding of 5th International Workshop on Information Hiding*, 2002.
  - [19] J. Fridrich, “Feature-based steganalysis for JPEG images and its implications for future design of steganographic schemes,” *Proceeding of 6th Information Hiding Workshop*, Toronto, 2004.
  - [20] H. Sajedi and M. Jamzad, “A steganalysis method based on contourlet transform coefficients,” *Proceeding of 4th International Conference on Intelligent Information Hiding and Multimedia Signal Processing*, 2008.
  - [21] J. Fridrich, T. Pevný, and J. Kodovský, “Statistically undetectable JPEG steganography: Dead ends, challenges, and opportunities,” *MM&Sec*, ACM, Dallas, USA, 2007.

# Research on DOA Estimation of Multi-Component LFM Signals Based on the FRFT

Haitao QU<sup>1</sup>, Rihua WANG<sup>2</sup>, Wu QU<sup>3</sup>, Peng ZHAO<sup>4</sup>

<sup>1</sup>Beijing University of Posts and Telecommunications, Beijing, China

<sup>2</sup>Communication University of China, Beijing, China

<sup>3</sup>K-Touch Corporation, Beijing, China

<sup>4</sup>Beijing Research Institute of China Telecom Co., Ltd, Beijing, China

E-mail: quhaitao2007@gmail.com

Received April 10, 2009; revised May 20, 2009; accepted May 25, 2009

## Abstract

A novel algorithm for the direction of arrival (DOA) estimation based on the fractional Fourier transform (FRFT) is proposed. Firstly, using the properties of FRFT and mask processing, Multi-component LFM signals are filtered and demodulated into a number of stationary single frequency signals. Then the one-dimensional (1-D) direction estimation of LFM signals can be achieved by combining with the traditional spectrum search method in the fractional Fourier (FRF) domain. As for the multi-component LFM signals, there is no cross-term interference, the mean square error (MSE) and Cramer-Rao bound (CRB) are also analyzed which perfects the method theoretically, simulation results are provided to show the validity of our method. The proposed algorithm is also extended to the uniform circular array (UCA), which realizes the two-dimensional (2-D) estimation. Using the characteristics of time-frequency rotation and demodulation of FRFT, the observed LFM signals are demodulated into a series of single frequency ones; secondly, operate the beam-space mapping to the single frequency signals in FRF domain, which UCA in array space is changed into the virtual uniform circular array (ULA) in mode space; finally, the DOA estimation can be realized by the traditional spectral estimation method. Compared with other method, the complex time-frequency cluster and the parameter matching computation are avoided; meanwhile enhances the estimation precision by a certain extent. The proposed algorithm can also be used in the multi-path and Doppler frequency shift complex channel, which expands its application scope. In a word, a demodulated DOA estimation algorithm is proposed and is applied to 1-D and 2-D angle estimation by dint of ULA and UCA respectively. The detailed theoretical analysis and adequate simulations are given to support our proposed algorithm, which enriches the theory of the FRFT.

**Keywords:** DOA Estimation, The Fractional Fourier Transform, UCA, ULA, LFM

## 1. Introduction

In various applications of array signal processing such as radar, sonar, communications, and seismology, there is a growing interest in estimating the DOA of LFM signals by dint of time-frequency analysis tools. G. Wang [1] proposed an iterative algorithm based on time-compensation, but the initial estimate is necessary. Using interpolation in the spatial time-frequency distribution matrices (STFD's) [2], Gershman [3] extended the signal subspace technique and estimated effectively DOA of

LFM signals, however Gershman's approach presences model biases in addition to time consuming. The above Wigner-Ville distribution (WVD) based methods consequently suffer from the disturbance of cross-terms in the presence of multi-component signals.

Using a new time-frequency analysis tool-FRFT, direction estimation of LFM signals has been proposed in Reference [4]. However, only maximal energy concentration point is selected as estimate data, easily interfered by surroundings. In this paper, a new FRFT based algorithm is proposed. Firstly, Observed signals are separated into a number of single components by adding an adap-

tive filter in the FRF domain. Secondly, the separated components are demodulated into stationary signals. Finally, the 1-D DOA of LFM signals can be estimated by the traditional spectrum search method. This algorithm digs two dimensional time and frequency information without the initial estimate, frequency focusing and parameter partnership. With the increasing of the Signal-to-Noise ratio (SNR), the MSE is quite closed to the CRB [5], for multi-component signals, cross-terms and non-linear optimize operation are also avoided.

For the UCA widely used in the third generation mobile communication system, the time-frequency characteristics of the FRFT are combined with the beamforming technology in FRF domain, an algorithm for the 2-D DOA estimation of the multi-component LFM signals is also proposed. Compared with other methods, the precision is enhanced by a certain extent. Simulation verifies the method to be effective in the multipath and Doppler frequency shift existed complex channels.

## 2. Background Knowledge of FRFT

### 2.1. Definition and Properties of FRFT

Recently the FRFT attracts more and more attention in the signal processing society, in 1980, Namias [6] firstly introduced the mathematical definition of the FRFT. Then Almeida [7] analyzed the relationship between the FRFT and the WVD, and interpreted it as a rotation operator in the time-frequency plane. This characteristic makes FRFT especially suitable for the processing of LFM signals [8–9].

As a generalization of the standard Fourier transform, the FRFT can be regarded as a counterclockwise rotation of the signal coordinates around the origin in the time-frequency plane. If the traditional Fourier transform of a signal can be considered as a  $\pi/2$  counterclockwise rotation from the time axis to the frequency axis, the FRFT can be accordingly considered as a counterclockwise rotation from the time axis to the  $u$  axis with an angle  $\alpha$ , as illustrated by Figure 1.

The FRFT of signal  $x(t)$  is represented as

$$X_\alpha(u) = F^p[x(t)] = \int_{-\infty}^{\infty} x(t) K_\alpha(t, u) dt \quad (1)$$

where  $p$  is called the order of the FRFT,  $\alpha = p\pi/2$ ,  $F^p[\bullet]$  denotes the FRFT operator and  $K_\alpha(t, u)$  is the kernel function of the FRFT

$$K_\alpha(t, u) = \begin{cases} \sqrt{\frac{1-j\cot\alpha}{2\pi}} \exp(j\frac{t^2+u^2}{2}\cot\alpha - jtu\csc\alpha), & \alpha \neq n\pi \\ \delta(t-u), & \alpha = 2n\pi \\ \delta(t+u), & \alpha = (2n+1)\pi \end{cases} \quad (2)$$

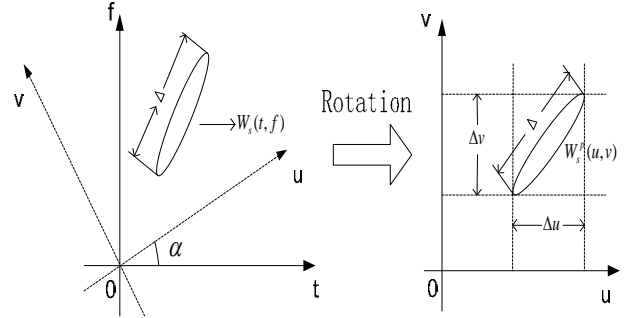


Figure 1. FRFT and WVD.

This has the following properties,

$$K_{-\alpha}(t, u) = K_\alpha^*(t, u) \quad (3)$$

$$\int_{-\infty}^{\infty} K_\alpha(t, u) K_\alpha^*(t, u') dt = \delta(u - u') \quad (4)$$

Hence, the inverse FRFT is

$$x(t) = F^{-p}[X(u)] = \int_{-\infty}^{\infty} X(u) K_{-\alpha}(t, u) du \quad (5)$$

Equation (5) indicates that signal  $x(t)$  can be interpreted as decomposition to a basis formed by the orthonormal LFM functions in the  $u$  domain, and the  $u$  domain is usually called the fractional Fourier domain, in which the time and frequency domains are its special cases. The FRFT is a one-dimension linear transform and has the rotation-addition property. Essentially, the representation of a signal in the fractional domains contains the information in both time and frequency domains of the signal; Thus the FRFT is considered as a time-frequency analysis method and has close relationships with other time-frequency analysis tools.

In Reference [10], some important characteristics are expressed as

$$F^p[e^{jct^2/2}] = \sqrt{\frac{1+j\tan\alpha}{1+c\tan\alpha}} \exp\left(\frac{u^2}{2} \frac{c-\tan\alpha}{1+c\tan\alpha}\right) \quad (6)$$

$$F^p[x(t)e^{jvt}] = X_p(u - v\sin\alpha) \sqrt{b^2 - 4ac} \exp\left[-j\left(\frac{v^2}{2}\sin\alpha\cos\alpha + uv\cos\alpha\right)\right] \quad (7)$$

$$F^p[x(t-\tau)] = X_p(u - \tau\cos\alpha) \exp\left[j(\tau^2\sin\alpha\cos\alpha/2 - u\tau\sin\alpha)\right] \quad (8)$$

### 2.2. Discrete FRFT Computation

In engineering applications, the discrete FRFT (DFRFT) is usually required. According to the definition of the FRFT, it is obvious that the numerical computation of the DFRFT is much more complicated than that of DFT. So far, there have been several DFRFT algorithms with

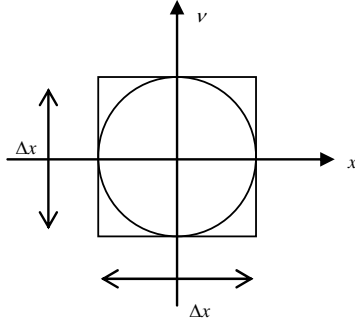


Figure 2. Normalized time-frequency support region.

different accuracies and different complexities. In this paper, we select the decomposition algorithm proposed in Reference [11]. This algorithm decomposes the computation of DFRFT to a convolution which can be computed by FFT, and the result is very close to the output of continuous FRFT. In this algorithm, the signal representation in time domain and frequency domain should be approximately constrained with an interval of  $[-T/2, T/2]$  and a bandwidth of  $[-F/2, F/2]$  respectively, viz. the time-bandwidth product of the signal is  $N = TF$ , and according to the uncertainty principle,  $N > 1$  constantly. If the sampling rate is selected as  $T_s = T/\sqrt{N}$ , the discrete representations of the signal in time domain and frequency domain will have the same length, which is called the dimensionless normalized process and the principle can be shown in Figure 2.

Therefore, Equation (1) can be expressed as

$$X_\alpha(u) = A_\alpha e^{j\pi u^2 \cot \alpha} \int_{-\infty}^{\infty} e^{j2\pi \csc \alpha ut} e^{j\pi t^2 \cot \alpha} x(t) dt \quad (9)$$

where

$$A_\alpha = \sqrt{\frac{1 - j \cot \alpha}{2\pi}}$$

For  $0.5 < |p| < 1.5$ , signal  $e^{j\pi t^2 \cot \alpha} x(t)$  has a bandwidth which is at most  $2F$  and can be represented using Shannon formula

$$e^{j\pi t^2 \cot \alpha} x(t) = \sum_{n=-N}^N e^{j\pi n^2 \cot \alpha / (2F)^2} x\left(\frac{n}{2F}\right) \cdot \text{sinc}\left(2F\left(t - \frac{n}{2F}\right)\right) \quad (10)$$

Substituting Equation (10) into Equation (9) and exchanging the sequence of the integral and the summation, we have

$$\begin{aligned} X_\alpha(u) &= F^p [x(t)] \\ &= \frac{A_\alpha}{2F} e^{j\pi u^2 \cot \alpha} \sum_{n=-N}^N e^{-j2\pi un \csc \alpha / (2F)} e^{j\pi n^2 \cot \alpha / (2F)^2} x\left(\frac{n}{2F}\right) \end{aligned} \quad (11)$$

By quantizing the variable  $u$  in the fractional Fourier domain, Equation (11) can be finally discretized as

$$X_\alpha(m) = F^p \left[ x\left(\frac{n}{2F}\right) \right] = \frac{A_\alpha}{2F} \sum_{n=-N}^N e^{j\pi(\gamma m^2 - 2\beta mn + \gamma n^2)/(2F)^2} x\left(\frac{n}{2F}\right) \quad (12)$$

where  $X_\alpha(m)$  denotes the DFRFT of signal  $x(t)$ ,  $\gamma = -\cot \alpha$ ,  $\beta = \csc \alpha$ . This algorithm can be implemented by FFT, and has a computation complexity of  $O(N \log_2 N)$  [11].

### 2.3. Two Special FRF Domain

WVD is an important non-stationary signal analysis tool, which has a very simple relationship with FRFT; viz. the WVD of FRFT is the coordinate rotation of the original signal' WVD, while the shape of WVD keeps unchanged in the rotation. Therefore, a lot of the WVD-based signal processing methods can be substituted by FRFT. The relationship of the two time-frequency analysis tools can draw a conclusion that "time width ( $\Delta u$ )" and "frequency width ( $\Delta v$ )" will change with the difference of the rotation angle. Considering two extreme cases,  $\Delta u \rightarrow 0, \Delta v \rightarrow \Delta$  or  $\Delta v \rightarrow 0, \Delta u \rightarrow \Delta$ , from the above analysis, the former corresponds to the rotation angle  $\alpha = -\cot^{-1} \mu$ , LFM signal becomes an impact function, which domain is called energy concentrated FRF one. The latter corresponds to the rotation angle  $\alpha = -\cot^{-1} \mu \pm \pi/2$ , LFM becomes a single frequency signal, which domain is called demodulated FRF one and is the base of the proposed algorithm in this paper. By dint of the time-frequency rotation property of FRFT, the detection, extraction and parameter estimation of LFM signals can be easily achieved.

### 2.4. The FRFT of Gaussian White Noise

**Theorem 1:** The FRFT of zero-mean Gaussian white noise is still Gaussian white noise.

**Proof:** let  $n(t)$  subject to the  $N(0, \sigma^2)$  distribution, and  $N_p(u)$  is its FRFT, the mean is

$$E\{N_p(u)\} = E\{F^p[n(t)]\} = F^p\{E[n(t)]\} = 0 \quad (13)$$

Because the FRFT is the linear transform, does not change the distribution characteristics of Gaussian noise. Therefore, the noise is still a zero mean Gaussian noise.

As for the second-order statistical properties of noise, the correlation of the white noise  $n(t)$  can be defined as:

$$E\{n(t)n^*(t-\tau)\} = \sigma^2 \delta(\tau) \quad (14)$$

The correlation of  $N_p(u)$  is defined as:

$$\begin{aligned}
 & E\{N_p(u)N_p^*(u-v)\} \\
 &= E\left\{\int_{-\infty}^{+\infty} n(t)K_p(t,u)dt \int_{-\infty}^{+\infty} n^*(\lambda)K_p^*(\lambda,u-v)d\lambda\right\} \\
 &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} E\{n(t)n^*(\lambda)\}K_p(t,u)K_p^*(\lambda,u-v)dt d\lambda \quad (15) \\
 &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \sigma^2 \delta(t-\lambda)K_p(t,u)K_p^*(\lambda,u-v)dt d\lambda \\
 &= \sigma^2 \int_{-\infty}^{+\infty} K_p(\lambda,u)K_p^*(\lambda,u-v)d\lambda
 \end{aligned}$$

Submit Equation (2) to Equation (15), and obtain:

$$\begin{aligned}
 & E\{N_p(u)N_p^*(u-v)\} \\
 &= \sigma^2 \left| \sqrt{\frac{1-j\cot\alpha}{2\pi}} \right|^2 e^{j0.5(2uv-v^2)\cot\alpha} \int_{-\infty}^{+\infty} e^{-j\lambda v \csc\alpha} d\lambda \quad (16) \\
 &= 2\pi\sigma^2 \left| \sqrt{\frac{1-j\cot\alpha}{2\pi}} \right|^2 |\sin\alpha| \delta(v)
 \end{aligned}$$

Due to Equation (16), we can see that the FRFT does not change the time-domain white characteristics of noise, while noise energy does not be changed.

Assume the array noise is the zero-mean airspace one, viz. as for the array element  $k$  ( $k \neq l$ ), the output noise is unrelated:

$$E\{n_k(t)n_l^*(t)\} = E\{n_k(t)\}E\{n_l^*(t)\} = 0 \quad (17)$$

The cross-correlation of the noise in FRF domain is

$$\begin{aligned}
 & E\{N_k^p(u)[N_l^p(u)]^*\} \\
 &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} E\{n_k(t)n_l^*(\lambda)\}K_p(t,u)K_p^*(\lambda,u)dt d\lambda \quad (18) \\
 &= 0
 \end{aligned}$$

The above equation shows, FRFT does not change the airspace white characteristics of noise. Therefore, we can draw a conclusion that FRFT does not change the statistical properties of Gaussian white noise, the theorem certification has completed.

Inference: as for the  $M$  antenna array element, if the array output noise is zero mean and variance  $\sigma^2$ , the noise covariance matrix in FRF domain is:

$$R_N^p = E\{N_p(u)N_p^*(u)\} = \sigma^2 I_M \quad (19)$$

### 3. 1-D DOA Estimation Algorithm

#### 3.1. ULA Array Model

Let a ULA of  $M$  sensors receive LFM sources from the  $D$  unknown directions  $\{\theta_1, \theta_2, \dots, \theta_D\}$ , as illustrated by

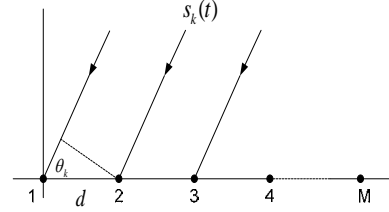


Figure 3. ULA and array model.

Figure 3. The observed signal at the output of the  $i$ th sensor can be described as

$$x_i(t) = \sum_{k=1}^D s_k[t - \tau_{ik}] + n_i(t) \quad (20)$$

$i = 1, 2, \dots, M-1 \quad k = 1, 2, \dots, D$

where,

$$s_k(t) = \exp[j(\omega_k t + \mu_k t^2 / 2)] \quad (21)$$

$$\tau_{ik} = (i-1)d \cos \theta_k / c \quad (22)$$

$\omega_k$ ,  $\mu_k$  are initial frequency and FM rate,  $s_k(t)$  is the  $k$ th source in reference sensor  $x_1$ .  $n_i(t)$  is the additive white Gauss noise with variance  $\sigma^2$ , which is assumed to be statistically independent with signal sources.  $\tau_{ik}$  is the  $k$ th's path delay,  $c$  is light velocity and  $d$  is sensor spacing.

From (20) and (21), we get the direction matrix is time-variant; however the traditional estimation method is merely suitable for time-invariant signal model. Therefore, the traditional method cannot be used to the direction finding of LFM signals directly.

#### 3.2. 1-D Estimation Algorithm Description

In this section, the main work is how to make the direction matrix time-invariant. The FRFT is actually a "Rotation" of signal in time-frequency plane. An LFM signal can be turned into an impulse in a proper fractional domain, for the ULA model, signal  $s_k(t)$  will present an impulse while the rotation angle  $\alpha_k = \pi - \cot \mu_k$ . There will be the energy concentration, consequently a distinct peak will appear in that FRF domain, whereas the noise energy is distributed much more symmetrically in the entire time-frequency plane and will not be concentrated in any FRF domain [12].

Using (20) and (21), we get that path delay can not change the FM rates, so the impulse corresponding rotation angles of signal  $s_k(t)$  are same in every sensor. Then Equation (20) is rotated with angle  $\alpha_k$  by the FRFT from two sides:

$$W_i^{\alpha_k}(u') = Y_{ik}^{\alpha_k}(u') + \sum_{l \neq k} Y_{il}^{\alpha_k}(u') + V_i^{\alpha_k}(u') \quad (23)$$



where,  $Y_{ik}^{\alpha_k}(u')$  presents an impulse,  $\sum_{l \neq k}^D Y_{il}^{\alpha_k}(u')$  and  $V_i^{\alpha_k}(u')$  are approximately considered as LFM signal and the white Gauss noise respectively.

Therefore, a mask operation is applied to (23) according to the peak position  $m_{ik}$ , which is a narrowband filter with central frequency  $m_{ik}$ , and with a properly selected bandwidth  $2L$ , most energy of the signal  $Y_{ik}^{\alpha_k}(u')$  will be removed. This procedure can be regarded as an open loop adaptive time-varying filter whose central frequency varies linearly following the peak position  $m_{ik}$ .

Signal  $Y_{ik}^{\alpha_k}(u')$  is performed the FFT (viz. FRFT of  $p=1$ ). According to the rotation-addition property [10], the two procedures above are equivalence to one time rotation with angle  $\alpha_k$  viz.

$$\begin{aligned}\alpha_k &= 3\pi/2 - \cot \mu_k; \\ \mu_k &= \tan \alpha_k\end{aligned}\quad (24)$$

Using (6), (7) and (8), signal  $s_k(t)$  is rotated with angle  $\alpha_k$  by the FRFT can be expressed as

$$\begin{aligned}S_k^{\alpha_k}(u) &= \frac{\sqrt{1+j\tan\alpha_k}}{\sqrt{1+\mu_k\tan\alpha_k}} \exp\left[\frac{(u-\omega_k\sin\alpha_k)^2}{2} \frac{\mu_k-\tan\alpha_k}{1+\mu_k\tan\alpha_k}\right] \\ &\quad \exp[-j(\omega_k^2\sin\alpha_k\cos\alpha_k/2 + u\omega_k\cos\alpha_k)] \\ &= \frac{\sqrt{1+j\tan\alpha_k}}{\sqrt{1+\mu_k\tan\alpha_k}} \exp[-j(\omega_k^2\sin\alpha_k\cos\alpha_k/2)] \\ &\quad \exp(-ju\omega_k\cos\alpha_k) = B \exp(-ju\omega_k\cos\alpha_k)\end{aligned}\quad (25)$$

where,

$$B = \frac{\sqrt{1+j\tan\alpha_k}}{\sqrt{1+\mu_k\tan\alpha_k}} \exp[-j(\omega_k^2\sin\alpha_k\cos\alpha_k/2)] \quad (26)$$

From (25) and (26), it can be seen that LFM signal  $s_k(t)$  has been transformed into the single frequency signal  $S_k^{\alpha_k}(u)$  in the FRF domain.

Similarly, the FRFT of path delayed signal  $s_k(t-\tau)$  with rotation angle  $\alpha_k$  can be expressed as

$$\begin{aligned}F^{\alpha_k}[s_k(t-\tau)] &= B \exp(j\tau^2\sin\alpha_k\cos\alpha_k/2) \\ &\quad \exp(j\tau\omega_k\cos^2\alpha_k) \exp[-ju(\omega_k\cos\alpha_k + \tau\sin\alpha_k)]\end{aligned}\quad (27)$$

In practice,  $\tau$  is too small viz.

$$\tau\sin\alpha_k \ll \omega_k\cos\alpha_k$$

$$\exp(j\tau^2\sin\alpha_k\cos\alpha_k/2) \approx 0 \quad (28)$$

Substituting (28) into (27), we get

$$\begin{aligned}F^{\alpha_k}[s_k(t-\tau)] &\approx B \exp(j\tau\omega_k\cos^2\alpha_k) \exp(-ju\omega_k\cos\alpha_k) \\ &= \exp(j\tau\omega_k\cos^2\alpha_k) S_k^{\alpha_k}(u)\end{aligned}\quad (29)$$

From the above analysis, Using (25) and (29), the observed signals described by (20) are performed the FRFT with rotation angle  $\alpha_k$  from two sides

$$\begin{aligned}X_{ik}^{\alpha_k}(u) &= S_{ik}^{\alpha_k}(u) + N_{ik}^{\alpha_k}(u) \\ i &= 1, 2, \dots, M-1\end{aligned}\quad (30)$$

Equation (30) can be compactly represented by matrix form as follows

$$X_k^{\alpha_k}(u) = A_k^{\alpha_k} S_k^{\alpha_k}(u) + N_k^{\alpha_k}(u) \quad (31)$$

$$A_k^{\alpha_k} = [a_{1k}, \dots, a_{ik}, \dots, a_{Mk}]^T \quad (32)$$

where,  $T$  denotes the transpose of matrix.

$$\begin{aligned}a_{ik} &= \exp(j\tau_{ik}\omega_k\cos^2\alpha_k) \\ &= \exp(j\frac{2\pi}{\lambda}\cos^2\alpha_k(i-1)d\cos\theta_k)\end{aligned}\quad (33)$$

$$X_k^{\alpha_k}(u) = [X_{1k}^{\alpha_k}(u), X_{2k}^{\alpha_k}(u), \dots, X_{Mk}^{\alpha_k}(u)]$$

$$N_k^{\alpha_k}(u) = [N_{1k}^{\alpha_k}(u), N_{2k}^{\alpha_k}(u), \dots, N_{Mk}^{\alpha_k}(u)] \quad (34)$$

From (32) and (33), the direction matrix  $A_k^{\alpha_k}$  is only relative to the direction information  $\theta_k$ , so the observed signal model has been time-variant in the FRF domain.

In the FRF domain, the covariance matrix of the observed signal can be defined as

$$R_{XX}^{\alpha_k} = E[X_k^{\alpha_k}(u)X_k^{\alpha_k H}(u)] = A_k^{\alpha_k} R_{SS}^{\alpha_k} A_k^{\alpha_k H} + \sigma^2 I \quad (35)$$

where,  $H$  denotes the conjugate transpose of matrix.  $R_{SS}^{\alpha_k}$  is the auto-correlation matrix of signal sources. The composite covariance matrix (35) has the same structure as the covariance matrix arising in the case of stationary signals. Therefore, the DOA can be estimated by performing eigendecomposition to  $R_{XX}^{\alpha_k}$ . Using the signal subspace  $S_N^{\alpha_k}$  and the noise subspace  $E_N^{\alpha_k}$ , the space spectrum function of the  $k$ th source in the FRF domain can be given by [13]

$$P(\theta_k) = 1 / (A_k^{\alpha_k H} E_N^{\alpha_k} E_N^{\alpha_k H} A_k^{\alpha_k}) \quad (36)$$

$P(\theta_k)$  is performed an 1-D search and  $\alpha_k$  can be obtain by the maximal peak rotation angle. Similarly, all the Direction of LFM signals can be estimated in turn. This

algorithm is considered as FRFT based demodulation method.

To summarize, the proposed algorithm can be formulated as follows:

- 1) The observed signals at all sensors are rotated with a continuously variable angle  $\alpha$  by the FRFT; perform a 2-D peak search in the  $(\alpha, m)$  plan to obtain the maximal peak position  $m_{ik}$  and corresponding rotation angle  $\alpha_k$  respectively.
- 2) Mask operations are applied according to  $m_{ik}$  at every sensor, then the filtered  $2L$  points are performed the FFT to obtain stationary signals consequently.
- 3) Get the covariance matrix of the stationary signals and perform eigendecomposition in the FRF domain, construct the spectrum function  $P(\theta_k)$  according to (36).
- 4) Perform 1-D peak search to  $P(\theta_k)$  and obtain the DOA of the  $k$ th LFM signal.
- 5) For multi-component LFM signals, all the direction can be estimated by repeating the above procedures.

## 4. 2-D DOA Estimation Algorithm Using UCA

### 4.1. Introduction

UCA has many advantages which the linear array cannot match. E.g. UCA can be implemented with all-direction-funding; its precision measurement does not change with the azimuth significantly and is fit for the system correcting. UCA is the main receiving antenna of base station system in the third generation mobile communication system. Thus, the UCA based DOA estimation has been a research hotspot in array signal processing. Mathwes [14] proposed an UCA-RB-MUSIC method, which can be only suitable for the stationary signals; however, the actually existed signals are non-stationary ones which are represented by LFM. Tao ran [4] proposed an algorithm of LFM signal DOA estimation. However, the method does not apply to the UCA.

Due to the above analysis, we propose a novel DOA estimation algorithm based on FRFT using UCA, as for the multi-component LFM signals, using the characteristics of time-frequency rotation and demodulation of FRFT. Firstly, the observed signals are demodulated into a series of single frequency ones; secondly, operate the beam-space mapping to the single frequency signals in FRF domain, which UCA in array space is changed into the virtual ULA in mode space; finally, the DOA estimation can be realized by the traditional spectral estimation method. The proposed algorithm mines the time, frequency and spatial information maximally; compared

with other method, the complex time-frequency cluster and the parameter matching computation are avoided; meanwhile enhance the precision [15]. As for the multi-component LFM signals, there is no cross-term interference, the proposed algorithm is also applicable for the multi-path and Doppler frequency shift channels.

### 4.2. UCA Array Model

Assuming  $D$  independent LFM signals and the pitch and azimuth angle is  $\{(\theta_1, \phi_1), (\theta_2, \phi_2), \dots, (\theta_D, \phi_D)\}$  respectively, the array element number of UCA is  $N$  and radius is  $r$ , the center is the reference point of receiving antenna, as shown in Figure 4. Then the output of the  $i$ th sensor is:

$$x_i(t) = \sum_{k=1}^D s_k[t - \tau_{ik}] + n_i(t) \quad (37)$$

$$i = 1, 2, \dots, N \quad k = 1, 2, \dots, D$$

where,

$$s_k(t) = \exp[j(\omega_k t + \mu_k t^2 / 2)] \quad (38)$$

$$\tau_{ik} = r \sin \theta_k \cos(\phi_k - \varepsilon_i) / c \quad (39)$$

$$\varepsilon_i = 2\pi(i-1) / N \quad (40)$$

$s_k(t)$  is the  $k$ th LFM source, and  $\omega_k$  and  $\mu_k$  are the initial frequency and FM rate respectively,  $\tau_{ik}$  is the path delay and  $c$  is the light velocity.  $n_i(t)$  is the additive white Gaussian noise with zero mean and variance  $\sigma^2$ , which is independent with signals.

From the Equations (37) and (38), the direction matrix of observed signals is time-varying in UCA, while the traditional DOA estimation algorithm is only suitable for the time-invariant model, which cannot be used to deal with LFM signal directly.

### 4.3. 2-D Estimation Algorithm Description

From Equations (26) and (28), operate the FRFT to Equation (37) with the rotation angle  $\alpha_k$  from two sides:

$$X_{ik}^{\alpha_k}(u) = S_{ik}^{\alpha_k}(u) + N_{ik}^{\alpha_k}(u) \quad (41)$$

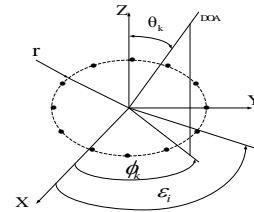


Figure 4. Uniform circular array.

The matrix form of Equation (41) is:

$$X^{\alpha_k}(u) = A^{\alpha_k} S_k^{\alpha_k}(u) + N^{\alpha_k}(u) \quad (42)$$

$$A^{\alpha_k} = [a_{1k}, \dots, a_{ik}, \dots, a_{Nk}]^T \quad (43)$$

where,  $T$  donates the transpose of matrix.

$$\begin{aligned} a_{ik} &= \exp(j\tau_{ik}\omega_k \cos^2 \alpha_k) \\ &= \exp(j2\pi r \sin \theta_k \cos(\phi_k - \varepsilon_i) \cos^2 \alpha_k / \lambda) \end{aligned} \quad (44)$$

From Equations (43) and (44), in appropriate FRF domain, direction matrix  $A^{\alpha_k}$  is only related to angle information  $\phi, \theta$ , viz. the observed signals have been transformed into unvaried smooth signal model. Therefore, the mode excitation method can be used to estimate the DOA of LFM signals.

The spatial beam former  $F_r^H$  in FRF domain is defined as

$$F_r^H = Q^H C e R^H \quad (45)$$

where,  $H$  denotes the conjugated transpose of matrix.

$$C e = \text{diag} \{ j^{-M}, \dots, j^{-1}, j^0, j^1, \dots, j^M \} \quad (46)$$

$$R^H = \sqrt{N} (Q_{-M}, \dots, Q_0, \dots, Q_M)^H \quad (47)$$

Select the central Hilbert matrix,

$$Q = \frac{1}{\sqrt{M'}} [v(\beta_{-M}), \dots, v(\beta_0), \dots, v(\beta_M)] \quad (48)$$

$$v(\psi) = [e^{-jM\psi}, \dots, e^{-j\psi}, e^{j0}, e^{j\psi}, \dots, e^{jM\psi}] \quad (49)$$

$$\beta_t = 2\pi / M' \quad t \in [-M, M] \quad (50)$$

where, the largest model number  $M \approx kr$ ,  $M' = 2M + 1$ . Wave number  $k = 2\pi / \lambda$ ,  $\lambda$  is the initial frequency corresponding center wavelength of LFM signals.  $F_r^H$  can change the UCA in the array space into the virtual ULA in the mode space, and finally, the DOA estimation can be achieved by the eigendecomposition based search method.

Summarize the above and the main steps are as follows:

- 1) The observed signals are continuously operated by FRFT; perform a 2-D peak search in the  $(\alpha, m)$  plan to obtain the maximal peak position  $m_{ik}$  and corresponding rotation angle  $\alpha_k$  of the  $k$  th LFM signal respectively.
- 2) Select  $2L$  points whose center is  $m_{ik}$  and calculate the FFT (FRFT with  $p = 1$ ), obtain the  $k$  th single frequency signal  $X_{ik}^{\alpha_k}(u)$ .
- 3) Let  $X_{ik}^{\alpha_k}(u)$  pass the beam switch  $F_r^H$ , viz.

$$\begin{aligned} Y_{ik}^{\alpha_k}(u) &= F_r^H X_{ik}^{\alpha_k}(u) \\ &= F_r^H A^{\alpha_k} S_k^{\alpha_k}(u) + F_r^H N^{\alpha_k}(u) \end{aligned} \quad ,$$

And calculate its covariance matrix

$$R_Y = E[Y_{ik}^{\alpha_k}(u) Y_{ik}^{\alpha_k}(u)^H] .$$

- 4) Define  $R = \text{Re}(R_Y)$ , perform eigendecomposition to  $R$  and obtain the signal subspace  $S$  and noise subspace  $G$ . Construct:

$$P(\theta_k, \phi_k) = \frac{1}{a_{ik}^T(\theta_k, \phi_k) G G^T a_{ik}(\theta_k, \phi_k)} ,$$

where,  $a_{ik}(\theta_k, \phi_k) = F_r^H a(\theta_k, \phi_k)$ , perform 2-D spectrum search and obtain  $\theta_k$  and  $\phi_k$ .

- 5) As for the multi-component LFM signal, repeat the above process and obtain all the DOA of signals respectively.

## 5. Performance Analysis and Simulation

### 5.1. FRFT Property Simulation

#### 5.1.1. Simulation of FRFT and WVD

As we all know, WVD is also one of the most important and most widely used time-frequency analysis tool, which is bound to FRFT with the existence of close ties. The derivation process is relatively complex; however, there is a very simple relationship between FRFT and WVD, that is, FRFT of WVD is the coordinate's rotation form of WVD of original signal [10].

In order to validate the relationship between the FRFT and WVD, experiments of compute simulations are given. We assume a wideband LFM signal  $s(t)$  with a length of 1024, which is modeled as: initial frequency and FM rates are  $\omega = 9\text{MHz}$ ,  $\mu = -0.7\text{MHz} / \mu\text{s}$ , sample frequency is  $f_s = 50\text{MHz}$ . The WVD of  $s(t)$  is shown in Figure 5 (a),  $s(t)$  is performed the FRFT by the rotation angle  $0.15\pi$  and get the transformed signal  $S_{0.15\pi}(u)$ . The WVD of the transformed signal  $S_{0.15\pi}(u)$  is shown in Figure 5(b). Compared the two figures, it can be found that the WVD of  $S_{0.15\pi}(u)$  is just the rotation of the WVD of  $s(t)$  by angle  $0.15\pi$ , meanwhile the figure shape is invariable. So the FRFT is testified a kind of rotation arithmetic operators in the time-frequency plane.

#### 5.1.2. Two Special FRF Domain Simulation

$s(t) = \exp[j(\omega t + \mu t^2 / 2)]$ , signal model is:  $\omega_1 = 9\text{MHz}$ ,  $\mu_1 = -1400000\text{MHz} / \text{s}$ . Sampling rate  $f_s = 50\text{MHz}$ , the number of snapshots is 1024. Perform continuous FRFT to signal and operate spectrum peak search, in the appropriate FRF domain,  $s(t)$  shows the property of energy

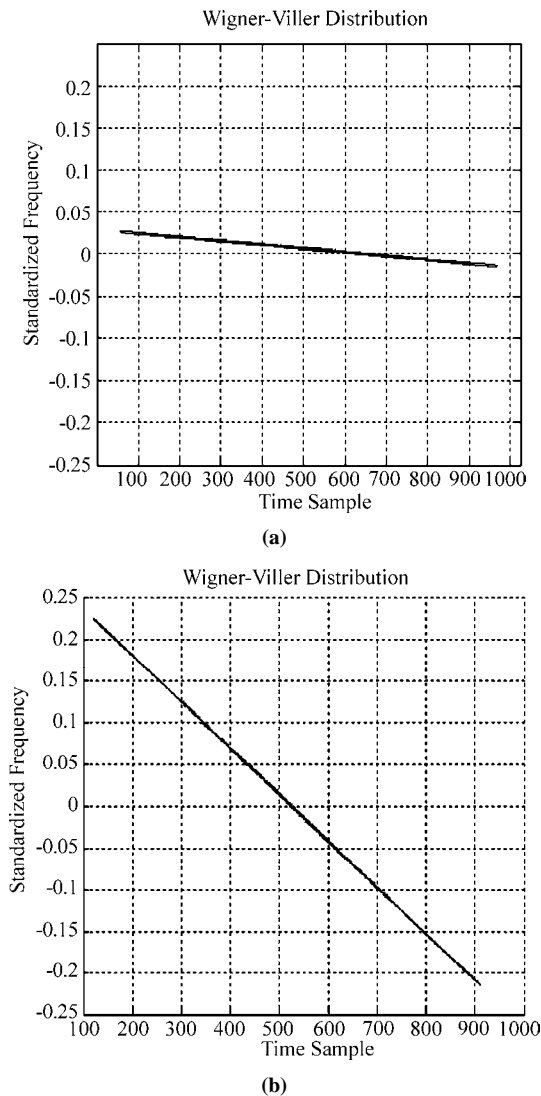


Figure 5. (a) The WVD of  $s(t)$ , (b) The WVD of  $S_{0.15\pi}(u)$ .

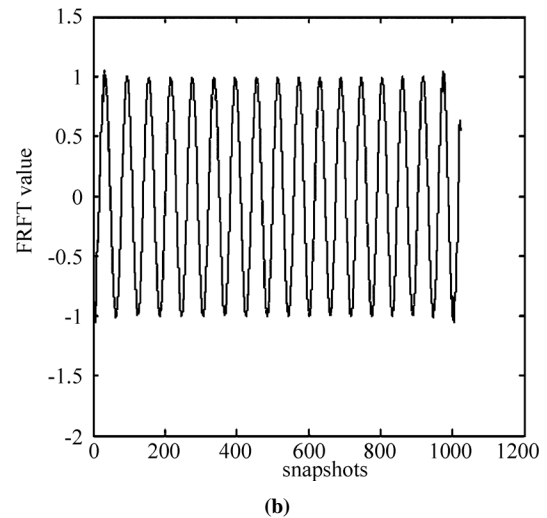
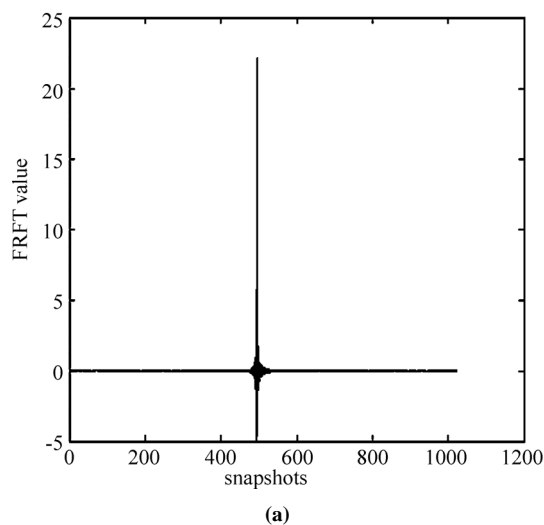


Figure 6. (a) Energy concentration property of FRFT, (b) Demodulated property of FRFT.

concentration, as shown in Figure 6(a). The signal continues to be rotated  $\pi/2$  in FRF domain, viz. in the demodulated FRF domain,  $s(t)$  shows the demodulated property, as shown in Figure 6(b).

### 5.1.3. Gaussian White Noise Simulation

Assume the complex Gaussian white noise is:  $w(n) = \text{randn}(1,1024) + j\text{randn}(1,1024)$  and perform continuous FRFT to it, the energy distribution of  $w(n)$  in different FRF domain is shown in Figure 7. We can see that the Gaussian white noise does not show energy concentration property in any FRF domain and can still be regarded as white noise. Thus theorem 1 is verified.

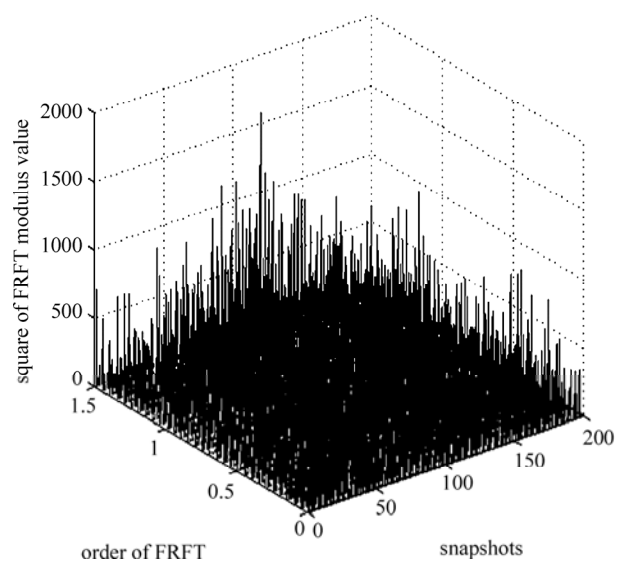


Figure 7. Energy distribution in different FRF domain.

## 5.2. 1-D DOA Estimation Simulation

### 5.2.1. MSE and CRB Analysis

The FRFT is a 1-D linear transform [10]. In the FRF domain  $N_k^{\alpha_k}(u)$  is approximately considered as the additive Gauss white noise. Therefore, the probability density function of signal  $X_k^{\alpha_k}(u)$  represents normal school and the corresponding likelihood function can be expressed as

$$L[X_k^{\alpha_k}] = \frac{1}{(2\pi)^M (\sigma^2/2)^M} \exp\left\{-\frac{1}{\sigma^2} [X_k^{\alpha_k} - A_k^{\alpha_k} S_k^{\alpha_k}]^H [X_k^{\alpha_k} - A_k^{\alpha_k} S_k^{\alpha_k}]\right\} \quad (51)$$

Using Reference [5], the CRB of the proposed method in the FRF domain can be represented as

$$CRB^{-1}(\theta_k) = \frac{2}{\sigma^2} \text{Re}\left\{(S_k^{\alpha_k})^H d_k^H(w) [I - A_k^{\alpha_k} (A_k^{\alpha_k H} A_k^{\alpha_k})^{-1} A_k^{\alpha_k H}] d_k(w) S_k^{\alpha_k}\right\} \quad (52)$$

where,  $\sigma^2$  is the noise variance and  $I$  is unit matrix,  $d_k(w) = dA_k^{\alpha_k} / dw$ .

Similarly, the MSE of the proposed algorithm in the FRF domain can be represented as

$$VAR_{MU}^{-1}(\theta_k) = \frac{2}{\sigma^2} [d_k^H(w) [I - A_k^{\alpha_k} (A_k^{\alpha_k H} A_k^{\alpha_k})^{-1} A_k^{\alpha_k H}] d_k(w) / \{[R_{XX}^{\alpha_k}]_{11} + \sigma^2 [R_{XX}^{\alpha_k} (A_k^{\alpha_k H} A_k^{\alpha_k})^{-1} R_{XX}^{\alpha_k}]_{11}\}] \quad (53)$$

where,  $R_{XX}^{\alpha_k}$  is covariance matrix of the observed signals,  $[ ]_{11}$  denotes the first row and first line element of matrix.

From (52) and (53), it can be obtain that the MSE of the proposed method will be more and more closed to the CRB with the increasing of the sensor number and the SNR.

### 5.2.2. MSE and CRB Simulation

In order to validate the proposed method, experiments of compute simulations are given. We assume the ULA of  $M=6$  impinging from  $D=2$  far field wideband LFM signals with a length of 1024, which is modeled as: initial frequency and FM rates are  $\omega_1 = 200\pi \text{ Hz}$ ,  $\mu_1 = -900\pi \text{ Hz/s}$ ;  $\omega_2 = 200\pi \text{ Hz}$ ,  $\mu_2 = 300\pi \text{ Hz/s}$ , angles of arrival are  $\theta_1 = 30^\circ$  and  $\theta_2 = 70^\circ$  respectively. Sample frequency is  $f_s = 900 \text{ Hz}$ , the mask snapshots are  $2L = 300$  in the FRF domain. The input SNR varies from 15dB to 29dB with an interval 2dB, at each level of the SNR, we run 100 Monte-Carlo experiments, the MSE of the proposed method and original method are

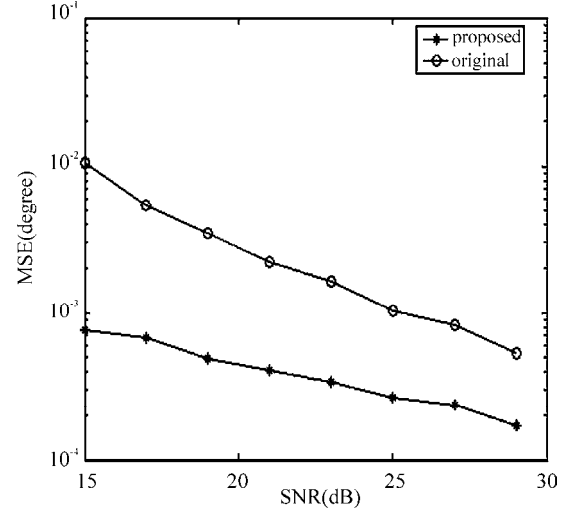


Figure 8. MSE of proposed and original method.

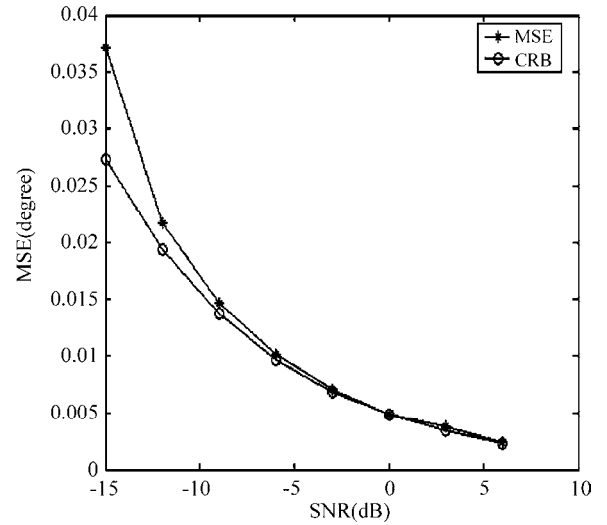


Figure 9. MSE and CRB of proposed method.

shown in Figure 8. Obviously, the accuracy of our method has certain improvement comparing with the method proposed in the Reference [4].

In same assumption, the input SNR varies from -15dB to 6dB with an interval 3dB, 100 times Monte-Carlo simulations are performed at each level of the SNR, MSE of the first signal and CRB are shown in Figure 9. It can be seen, the MSE of proposed method is closed to the CRB even at the lower SNR.

## 5.3. 2-D DOA Estimation Simulation

### 5.3.1. 2-D Estimation RMSE Simulation

$D=2$  two far-field LFM sources shoot the  $N=20$  UCA with the angle information  $\{(\theta_1 = 60^\circ, \beta_1 = 50^\circ), (\theta_1 = 30^\circ, \beta_1 = 70^\circ)\}$ . The signal model is:  $\omega_1 =$

$200\pi\text{Hz}$  ,  $\mu_1 = 300\pi\text{Hz}/s$  ;  $\omega_2 = 200\pi\text{Hz}$  ,  $\mu_2 = -900\pi\text{Hz}/s$  . Sampling rate is  $f_s = 900\text{Hz}$  , number of snapshots is 1024, and the cover filter length is  $2L = 300$  . The Figure 10(a) gives the 2-D DOA estimation of signal one in the 0dB SNR.

Change the input SNR range from 0dB to 20dB with the interval 5dB, firstly perform big step search to obtain the rough DOA estimation. Then run the high differentiation search with the 0.001rad step. Run 300 time Monte-Carlo experiment respectively, the RMSE (root mean square error, RMSE) comparison curves of the proposed algorithm and literature one can be seen in Figure 10(b). The accuracy of our method has certain improvement compared to the original algorithm.

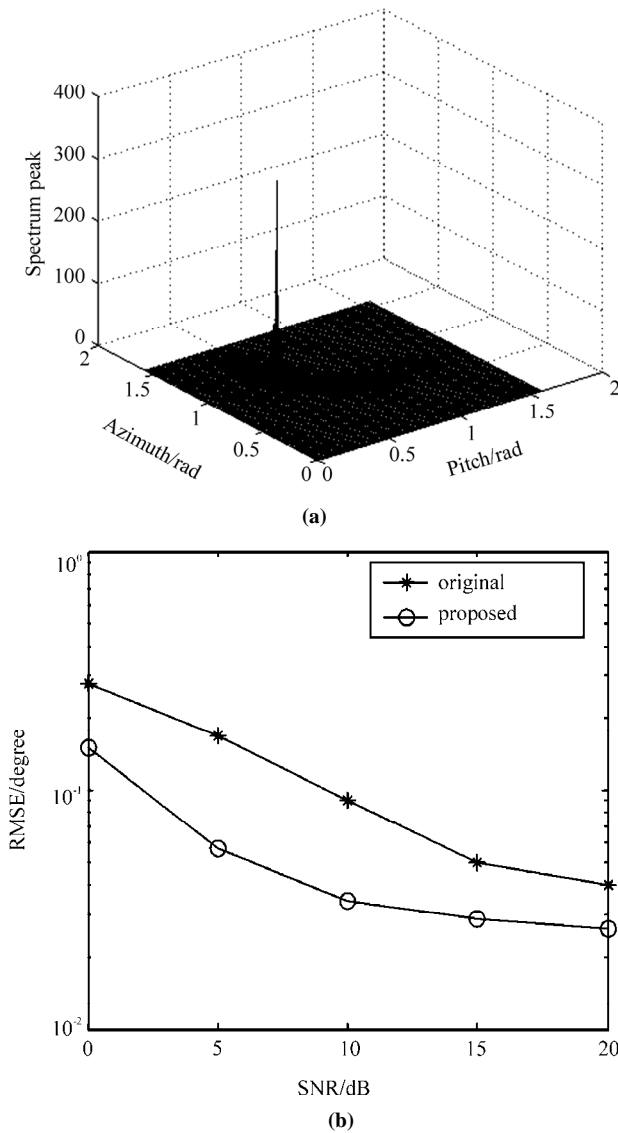


Figure 10. (a) 2-D DOA estimation using UCA, (b) RMSE comparison curves using UCA.

### 5.3.2. 2-D Estimation Performances in Complex Channel and Simulation

In mobile communication system, the proposed algorithm is applied to the complex channel which the multi-path and Doppler shift is existed simultaneously. In the same simulation conditions, viz. the random signal source model is:

$$s_k(t) = \sum_{e=1}^E M_e \exp(jf_e t) \exp[j(\omega_k(t - \delta_e) + \mu_k(t - \delta_e)^2 / 2)] \quad (54)$$

where,  $M_1 = 1, M_2 = 0.9$  , Doppler frequency shift is  $f_1 = 0, f_2 = 2$  , multi-path delay is  $\delta_1 = 0, \delta_2 = 1/900$  . When the SNR is 0dB, the simulation result of signal one in most powerful path can be shown in Figure 11(a).

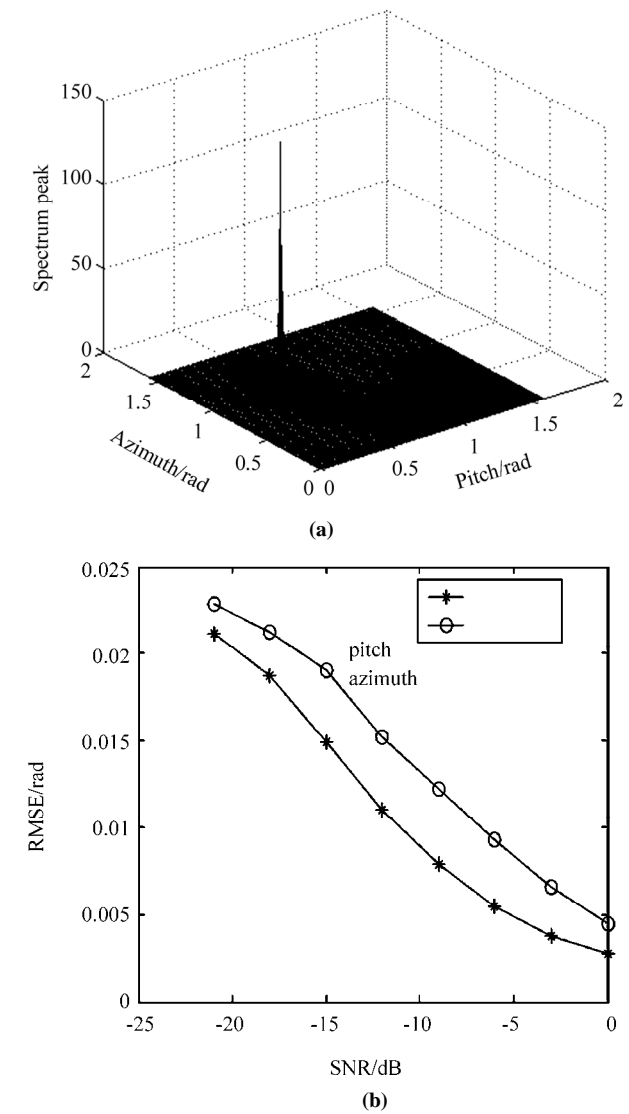


Figure 11. (a) 2-D DOA estimation in complex channel, (b) RMSE curves in complex channel.

Change the input SNR range from  $-21\text{dB}$  to  $0\text{dB}$  with the interval  $3\text{dB}$ . Run 300 time Monte-Carlo experiment respectively, the RMSE comparison curves of signal one can be seen in Figure 11(b), which can show that the proposed algorithm is also effective in complex channel.

## 6. Conclusions

Analyzing the definition and characteristics of the FRFT, a novel DOA estimation algorithm has been presented. In the implementation of the method, mask operation is introduced to simplify the filtering procedure with no accuracy degradation. Demodulation operation is used to extend the application range of the traditional estimate method without performance loss. Compared with other methods, the veracity has certain improvement while the cross-terms and interpolation are avoided. The proposed is also expanded to the 2-D DOA estimation using UCA, which is suitable for the multi-path and Doppler frequency shift complex environment. In addition, the proposed method can be also applied to DOA estimation of LFM signals in colored noise or near-field environment, which is not described in this paper.

The theoretical analyses about the error and CRB are also provided and verified by simulation results thus making this method more reliable in theory and in practice, meanwhile enrich the principle and application of the FRFT. As for the optimization, the 2-D Cramer-Rao Bound of the proposed algorithm is the further research direction.

## 7. Acknowledgements

The authors would like to thank the reviewers for their detailed comments on earlier versions of this paper.

## 8. References

- [1] G. Wang and X. G. Xia, "Iterative algorithm for direction of arrival estimation with wideband chirp signals," *IEEE Proceedings of Radar, Sonar, Navig.*, Vol. 147, No. 5, pp. 233–238, 2000.
- [2] A. Belouchirani and M. G. Amin, "Time-frequency MUSIC [J]," *IEEE Signal Processing Letters*, Vol. 6, No. 5, pp. 109–110, 1999.
- [3] A. B. Gershman and M. G. Amin, "Wideband direction of multiple chirp signals using spatial time-frequency distributions," *IEEE Signal Processing Letters*, Vol. 7, pp. 152–155, June 2000.
- [4] R. Tao and Y. S. Zhou, "A novel method for the DOA estimation of wideband LFM sources based on FRFT," *Transactions of Beijing Institute of Technology*, Vol. 25, No. 10, pp. 895–899, 2005.
- [5] P. Stoica and A. Nehorai, "Music, maximum likelihood, and cramer-rao bound," *IEEE Transactions on ASSP*, Vol. 37, No. 5, May 1989.
- [6] V. Namias, "The fractional Fourier transform and its application in quantum mechanics [J]," *IMA Journal of Applied Mathematics*, No. 25, pp. 241–265, 1980.
- [7] L. B. Almeida, "Fractional Fourier transform and time-frequency representations [J]," *IEEE Transactions on Signal Processing*, Vol. 42, No. 11, pp. 3084–3091, 1994.
- [8] Y. Q. Dong, R. Tao, S. Y. Zhou, *et al.*, "SAR moving target detection and imaging based on fractional Fourier transform," *Acta Armamentarii* (in Chinese), Vol. 20, No. 2, pp. 132–136, 1999.
- [9] L. Qi, R. Tao, S. Y. Zhou, *et al.*, "Adaptive time-varying filter for linear FM signal in fractional Fourier domain," *Proceedings of the 6th ICSP*, Posts and Telecommunications Press, Beijing, pp. 1425–1428, 2002.
- [10] R. Tao, L. Qi, and Y. Wang, "Principle and application of the fractional Fourier transform," Tsinghua Publishing Company, Beijing, 2004.
- [11] H. M. Ozaktas, O. Arikan, M. A. Kutay, *et al.*, "Digital computation of the fractional Fourier transform," *IEEE Transactions on Signal Processing*, Vol. 44, No. 9, pp. 2141–2150, 1996.
- [12] L. Qi, R. Tao, S. Y. Zhou, *et al.*, "Detection and parameter estimation of multicomponent LFM signal based on the fraction Fourier transform [JJ]," *Science in China (Series E)*, Vol. 47, No. 2, pp. 184–198, 2004.
- [13] X. D. Zhang, *et al.*, *Modern Signal Processing* (Second editor), Tsinghua Publishing Company, Beijing, 2002.
- [14] C. P. Mathews, "Eigenstructure techniques for 2-D angle estimation with uniform circular arrays," *IEEE Transactions on Signal Processing*, Vol. 42, No. 9, pp. 2395–2407, September 1994.
- [15] L. M. Yang "DOA estimation for wideband sources based on UCA," *Journal of Electronics, China*, Vol. 23, No. 1, January 2006.

# Novel Rate-Control Algorithm Based on TM5 Framework

Zhongjie ZHU<sup>1,2</sup>, Yongqiang BAI<sup>1,2</sup>, Zhiyong DUAN<sup>1,2</sup>, Feng LIANG<sup>1</sup>

<sup>1</sup>Ningbo Key Laboratory of DSP, Zhejiang Wanli University, Ningbo, China

<sup>2</sup>Physical Engineering College, Zhengzhou University, Zhengzhou, China

E-mail: zhongjiezhu@hotmail.com

Received April 9, 2009; revised April 25, 2009; accepted May 30, 2009

## Abstract

Rate control is a key technology in the fields of video coding and transmission, and it has attracted a great attention and has been studied extensively. The TM5 framework of MPEG-2 is a classical rate control algorithm and has been widely used. However, it has some underlying drawbacks during practical applications such as the poor rate control precision and high computational complexity. Hence, in this paper, a novel rate-control algorithm based on the TM5 framework is proposed. The drawback of the target bit allocation method of the original TM5 algorithm is firstly analyzed and improved. Then, a new rate-distortion model is incorporated into the rate control algorithm to implement rate prediction to enhance the rate-control precision. Meanwhile, the macro-block (MB) level rate control is adapted to be frame level to reduce the computational complexity. Experiments are conducted and some results are given. Compared with the original TM5 algorithm, the improved novel algorithm not only can enhance the rate-control precision but also can reduce the complexity and the fluctuation of decoded image quality.

**Keywords:** Rate Control, TM5, Target Bit Allocation, Rate-Distortion Model

## 1. Introduction

Rate control is a key technology in the fields of video coding and transmission [1,2]. With the rapid progress of video coding technology and explosion of video applications, it has attracted a great attention and has been studied extensively. The main objective of rate control is to optimally allocate available bits within video sequences to minimize visual distortion under the bit rate constraint. For a rate control algorithm, the rate-distortion performance and the computational complexity are two main issues that should be addressed. Although rate control is a normative part in video coding standards, almost all the main existing video coding standards have proposed their own recommendations on rate control over the last few years such as the TM5 algorithm of MPEG-2, the VM8 algorithm of MPEG-4, the TMN8 algorithm of H.263 and the F086/G012 of H.264 [3,4].

As mentioned above, the rate control is an informative part in video coding standard, which means that this part is still open for research. It leaves the flexibility for designers to develop suitable scheme for specific applications. Hence, this topic is still being studied extensively. Xu *et al* have proposed a novel Dynamic Video Rate

Control (DVRC) technique which can enable adaptive video delivery over the Internet [5]. But its performance in heterogeneous network environments should be further enhanced. Papadimitriou *et al* have proposed a novel rate control algorithm when hierarchical B-picture coding is used in H.264/AVC, where significant PSNR gains as well as accurate rate control precision can be achieved [6]. However, the computational load is high. For more other related works, the readers are referred to [2].

It is noted that, the TM5 rate control algorithm of MPEG-2 has obtained great attention and has been widely used. It implements rate control in macro-block level and mainly consists of three steps: target bit allocation, rate control and adaptive quantization. However, it also has drawbacks including poor rate-control precision and not low computational complex due to its macro-level quantification [7,8]. In this paper, a novel rate-control algorithm based on TM5 framework is proposed. The target bit-allocation of TM5 is improved and a new rate-distortion model is incorporated to implement rate prediction to enhance the rate-control precision. Meanwhile, the macro-block level rate control is adapted to frame level to reduce the computational complexity.

The remaining of the paper is organized as follows:



The TM5 rate control algorithm is briefly reviewed in Section 2; the target bit allocation is analyzed and an improved rate-distortion model is introduced in Section 3; the proposed novel algorithm is introduced in Section 4; Section 5 shows the experimental results to evaluate our work; Section 6 concludes the paper.

## 2. TM5 Rate Control Algorithm

The TM5 rate control algorithm has been designed for MPEG-2 standard. It mainly consists of the following three steps:

### 2.1. Target Bit Allocation

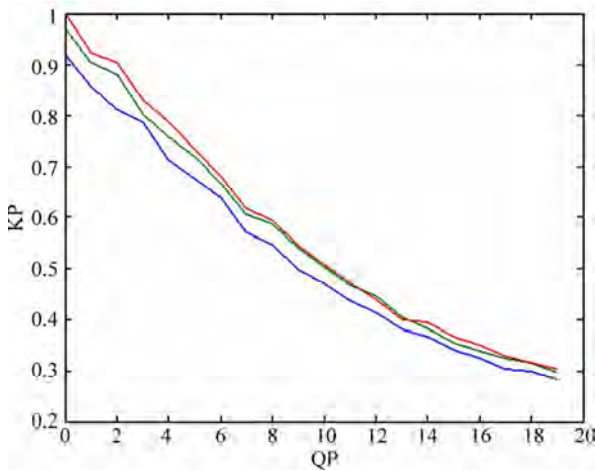
The target number of bits for the next picture depends on picture-type and “universal” weighting factors. The target number of bits for different type of frames ( $T_I$ ,  $T_P$ ,  $T_B$ ) are calculated by [9]:

$$T_I = \frac{R}{1 + \frac{N_P X_P}{K_P X_I} + \frac{N_B X_B}{K_B X_I}} \quad (1)$$

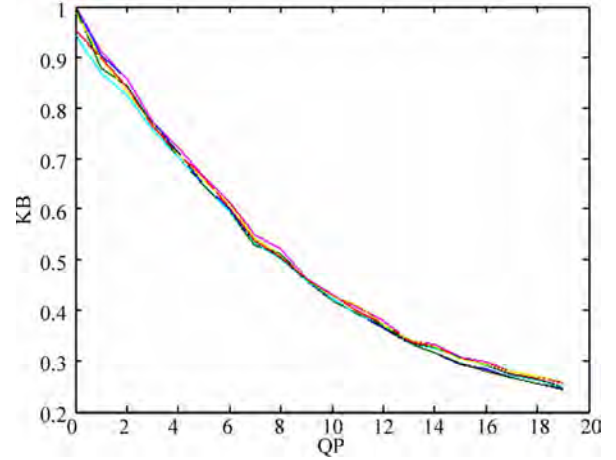
$$T_P = \frac{R}{N_P + \frac{N_B K_P X_B}{K_B X_P}} \quad (2)$$

$$T_B = \frac{R}{N_B + \frac{N_P K_B X_B}{K_P X_B}} \quad (3)$$

where  $R$  is the remaining number of bits assigned to GOP,  $N_P$ ,  $N_B$  are the number of P-pictures and B-pictures remaining in the current GOP,  $K_P$ ,  $K_B$  are universal constants depending on the quantization matrices. In most cases,  $K_P = 1.0$  and  $K_B = 1.4$ .

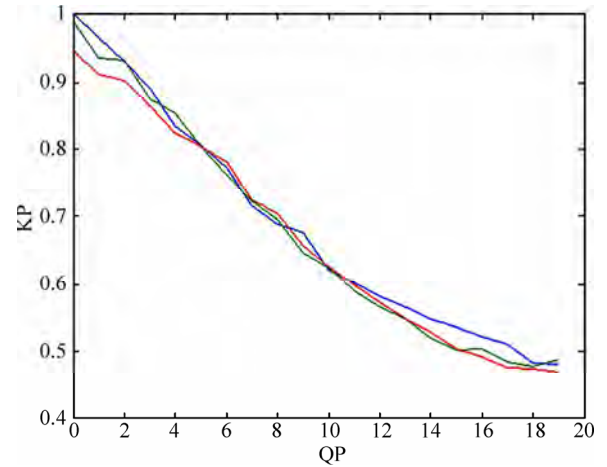


(a)  $K_P$ -QP curves of three P-frames

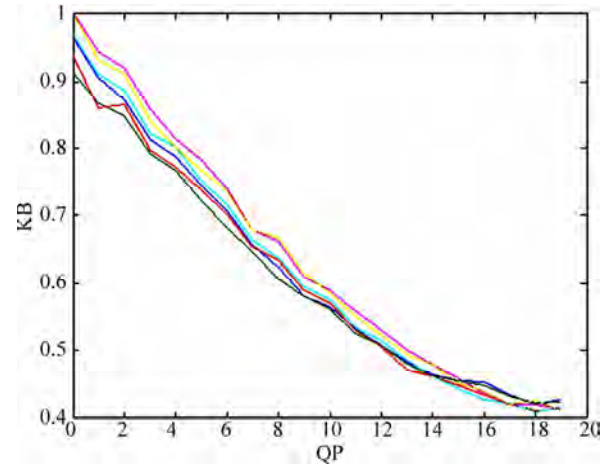


(b)  $K_B$ -QP curves of the six B-frames

Figure 1. The ( $K_P$ ,  $K_B$ )-QP curves of *Alex* sequence from one GOP.

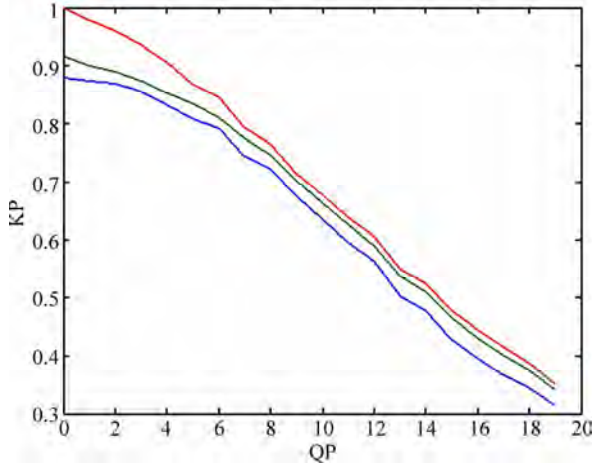
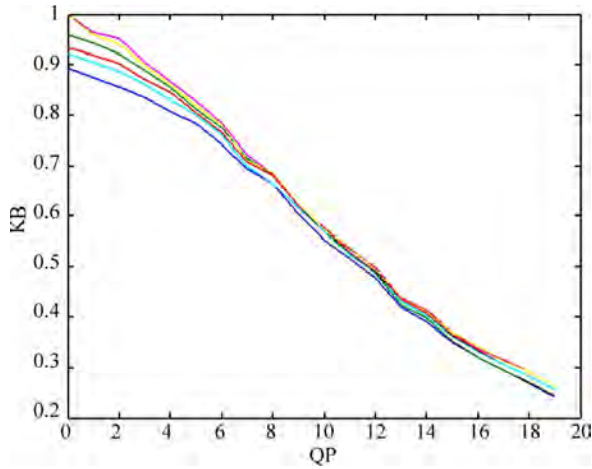


(a)  $K_P$ -QP curves of the three P-frames



(b)  $K_B$ -QP curves of the six B-frames

Figure 2. The ( $K_P$ ,  $K_B$ )-QP curves of the *Claire* sequence from one GOP.

(a)  $K_P$ -QP curves of the three P-frames(b)  $K_B$ -QP curves of the six B-frames

**Figure 3.** The  $(K_P, K_B)$ -QP curves of the *Train* sequence from one GOP.

## 2.2. Rate Control

The reference value of the quantization parameter for each macro-block (MB)  $Q_j$  is set as follows:

$$Q_j = \left( \frac{d_j \times 31}{r} \right) \quad (4)$$

where  $r$  is the reaction parameter and given by

$$r = 2 \frac{R}{f} \quad (5)$$

where  $R$  denotes the bit rate,  $f$  denotes the frame rate.

## 2.3. Adaptive Quantization

The final quantization parameter  $mquant_j$  for the  $j$ th

Macro-block  $MB_j$  is given by

$$mquant_j = Q_j \times N_{act_j} \quad (6)$$

where  $N_{act_j}$  is the normalized spatial activity measured for  $MB_j$ .

## 3. Improved Target Bit Allocation and Rate Distortion Model

### 3.1. Improved Target Bit Allocation

From Formulas (1), (2) and (3), it can be seen that  $K_P$ ,  $K_B$  are universal constants depending on the quantification matrix, which can be viewed as the ratio of number of bits of I frame to that of P frame or B frame, respectively, that is

$$K_P = \frac{R_I}{R_P}, K_B = \frac{R_I}{R_B} \quad (7)$$

where  $R_I$ ,  $R_P$ ,  $R_B$  denote the number of bits of I, P, B frames respectively.

In practical video coding applications,  $K_P$  and  $K_B$  are not constant, they are usually related to the quantification parameter of I-frame. Hence, rate control error and video quality fluctuation may increase if  $K_P$  and  $K_B$  keep constant during the whole encoding process.

To investigate the actual relations between  $K_P$ ,  $K_B$  quantification parameter of I-frame, experiments are conducted. Some standard test sequences, including Alex, Claire, and Train, are used to implement coding experiments to investigate this phenomenon. The GOP structure is set to IBBPBBPBBP [10]. Some experimental results are given in Figures 1–3 where the values of  $K_P$  and  $K_B$  are normalized and the actual  $K_P$ -QP curves of three P-frames and the  $K_B$ -QP curves of six B-frames from the same GOP are given. From the experimental results, it can be observed that the decrease of  $K_P$  or  $K_B$  is approximately linear to the increment of QP of I-frame.

Hence, the  $K_P$ -QP or  $K_B$ -QP relations can be approximately modeled as

$$K_i = a \times QP + b \quad i \in \{P, B\} \quad (8)$$

where  $K_i$  denotes the frame-level bit-allocation coefficient of P or B frames, QP is the quantification parameter of I frame,  $a$  and  $b$  are model parameters.

### 3.2. Novel Rate-Distortion Model

R-D models have been introduced since MPEG-4 and H.263. These models are helpful in rate control since they can provide sufficient information for determining quantification parameters. Once the target bit rate for the current encoding frame is acquired, the quantization parameter can be determined through the R-Q models. In order to improve the precision of the traditional quadratic R-D model, based on empirical observations and lots of experiments, we improved the quadric R-Q model and proposed a novel one [11]:

$$R = \frac{a}{\sqrt{Q}} + \frac{b}{Q^2} + c \quad (9)$$

where  $a$ ,  $b$ ,  $c$  are model parameters which can be calculated by linear regression method explained as follows. Define

$$x_1(Q) = \frac{1}{Q^2}, \quad x_2(Q) = \frac{1}{\sqrt{Q}} \quad (10)$$

Suppose  $(x_{11}, x_{21}, R_1)$ ,  $(x_{12}, x_{22}, R_2)$ , ...,  $(x_{1n}, x_{2n}, R_n)$  are existing samples, and define

$$M = \begin{pmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} \end{pmatrix}, R = \begin{pmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{pmatrix}, C = \begin{pmatrix} c \\ b \\ a \end{pmatrix} \quad (11)$$

Then based on linear regression method, the parameters can be estimated:

$$C = (M^T M)^{-1} M^T R, \quad (12)$$

where  $M^T$  is the transpose of  $M$  and  $(M^T M)^{-1}$  is the inverse matrix of  $M^T M$ .

## 4. Proposed Algorithm

Based on the above observations, the traditional TM5 rate control framework is adapted and a new rate control algorithm is proposed. The target bit-allocation is improved and a new rate-distortion model is incorporated to implement rate prediction to enhance the rate-control precision. Meanwhile, the macro-block level rate control is adapted to frame level to reduce the computational complexity. The key steps of the rate control are briefly introduced as follows:

### 4.1. Target Bits Calculation

Suppose  $N_G$  is the length of GOP,  $f_r$  is the frame-

rate,  $b_r$  is the bit-rate, the initial target number of bits for a GOP is calculated as:

$$T(0) = \frac{b_r}{2f_r} N_G - (B_0 - B_c(0)) \quad (13)$$

where  $B_0$  denotes the initial value of virtual buffer,  $B_c(j)$  denotes the occupancy of virtual buffer after encoding the  $j$ th frame,  $B_c(0)$  denotes the occupancy of virtual buffer after encoding the former GOP,  $T(j)$  denotes the remaining bits available for the current GOP after encoding the  $j$ th frame. After encoding one frame,  $T(j)$  is updated as:

$$T(j) = T(j-1) - A(j-1) \quad (14)$$

where  $A(j)$  is the number of bits generated by encoding the  $j$ th frame.

For frame-level bit-allocation, it is very important to properly select  $K_P$  and  $K_B$ . The bigger the values of  $K_P$  and  $K_B$  are, the smaller the distortion of encoded I frame will be. However, if  $K_P$  and  $K_B$  are set too large values, not only will the stream fluctuation increase but also video quality will decrease. Hence, in our algorithm,  $K_P$  and  $K_B$  are selected according to Equation (8).

### 4.2. Frame-Level Rate Control

After the target number of bits for each frame is estimated, the quantification parameter can be calculated through the following rate-distortion model:

$$\frac{R-H}{X} = \frac{a}{\sqrt{Q}} + \frac{b}{Q^2} + c \quad (15)$$

where  $a$ ,  $b$ ,  $c$  are model parameters,  $H$  is the head information. The complexity  $X$  of current frame is expressed by  $SAD$  predicted from that of the former frame, which is calculated as:

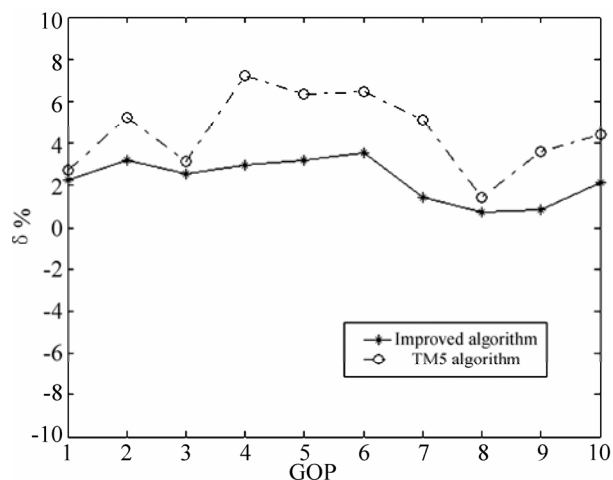
$$SAD = \sum_{(x,y)} \text{abs}(f(x,y) - \overline{f(x,y)}) \quad (16)$$

### 4.3. Model Update

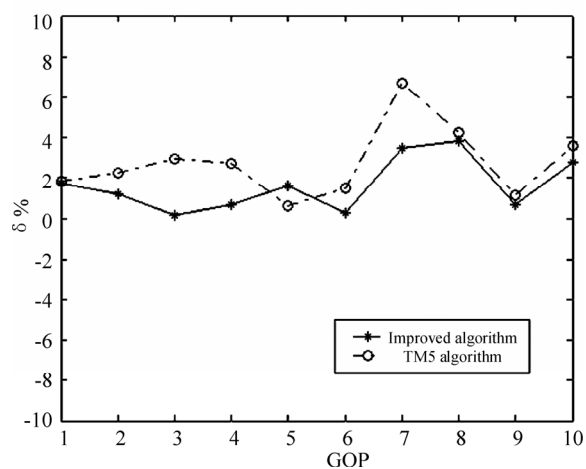
When finishing encoding one frame, the  $SAD$  model and the  $R-Q$  model are updated until the whole video sequence is encoded.

## 5. Experiment Results

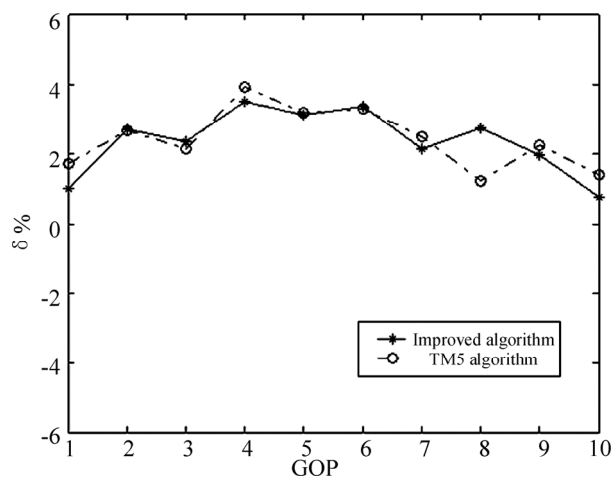
In order to evaluate the performance of the proposed rate control algorithm, we conduct theoretical analysis as



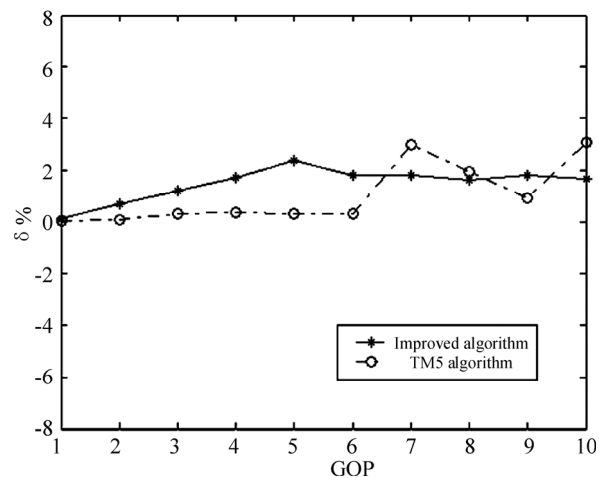
(a) 160Kbps



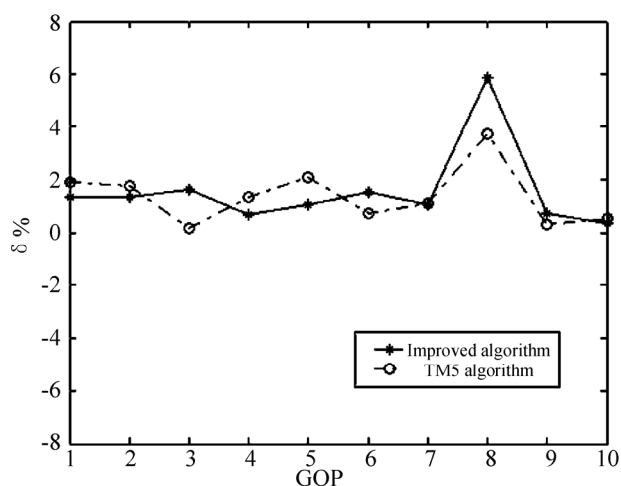
(a) 1.2Mbps



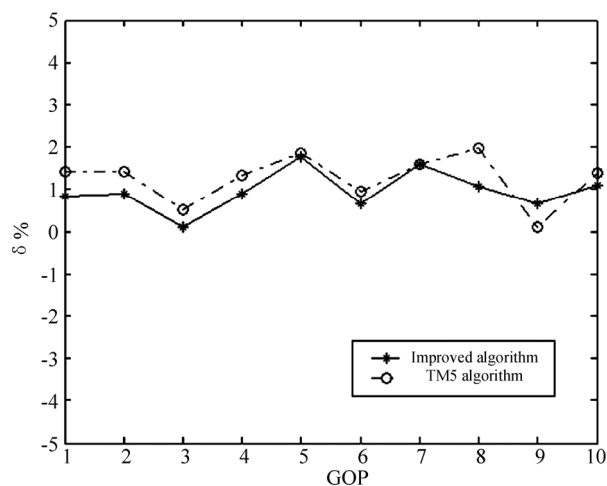
(b) 150Kbps



(b) 1.0Mbps



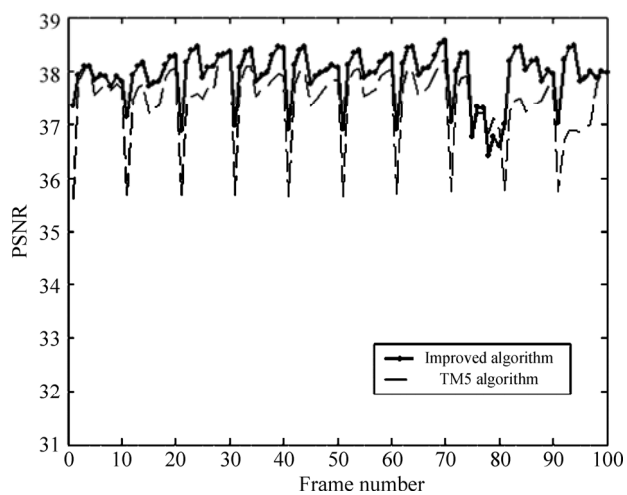
(c) 130Kbps



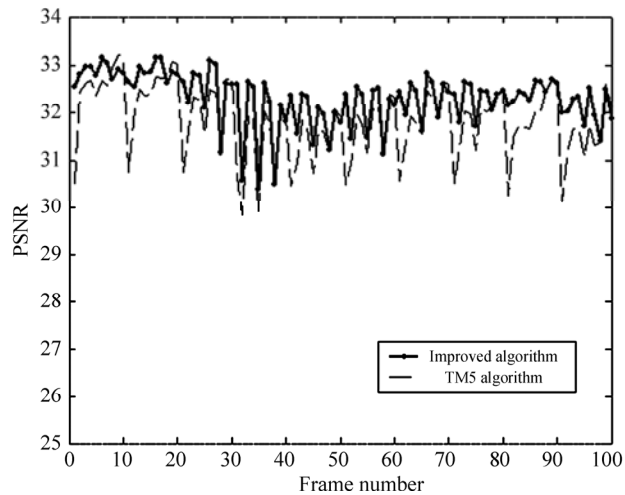
(c) 0.6Mbps

Figure 4. Rate-control precisions of the *Alex* sequence at different bit rates.

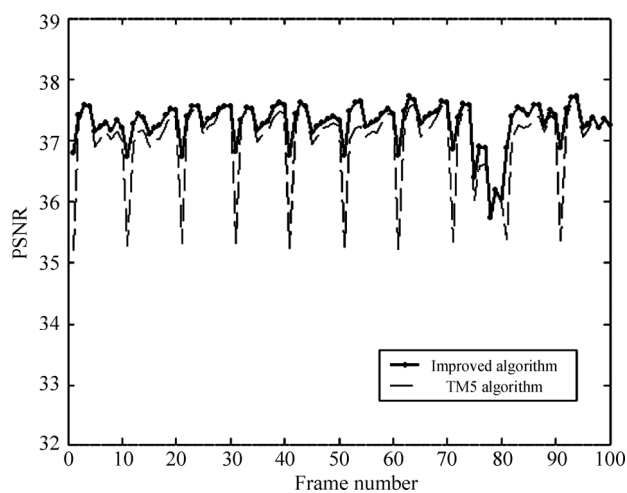
Figure 5. Rate-control precisions of the *Train* sequence at different bit rates.



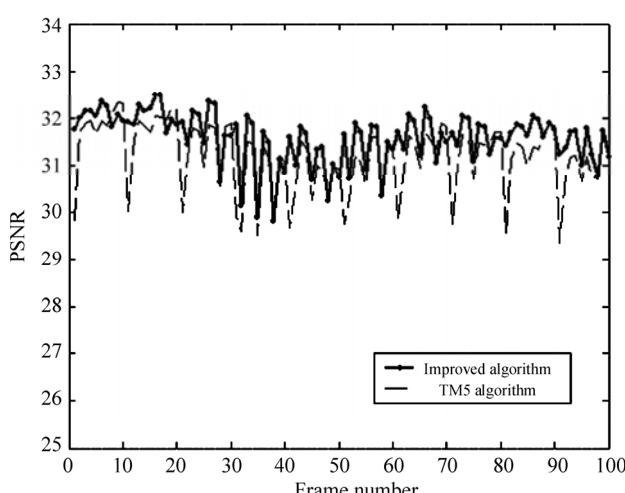
(a) 160Kbps



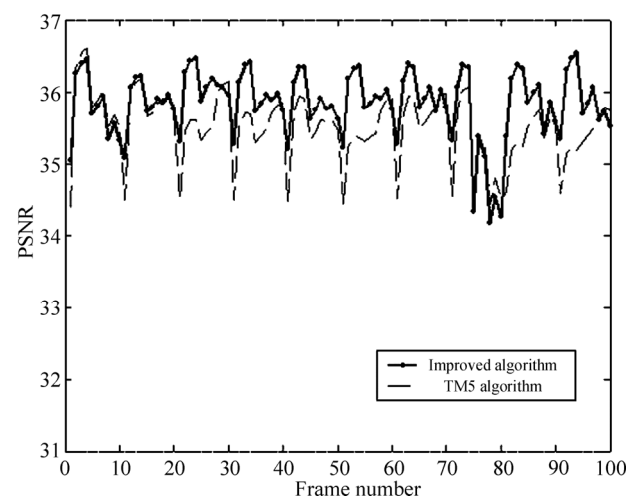
(a) 1.2Mbps



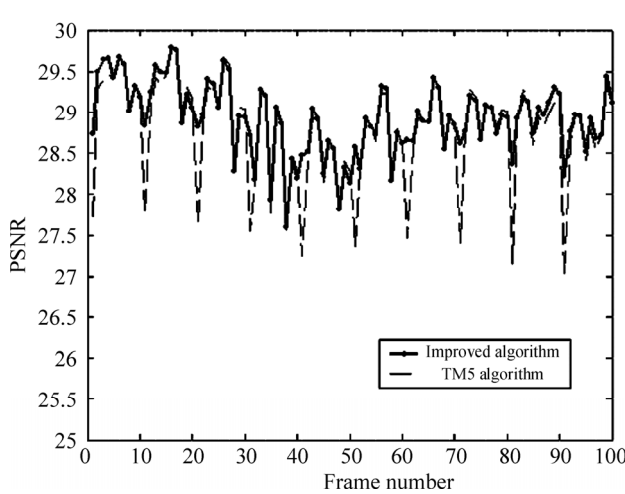
(b) 150Kbps



(b) 1.0Mbps



(c) 130Kbps



(c) 0.6Mbps

Figure 6. Image qualities of the *Alex* sequence at different bit rates.

Figure 7. Image qualities of the *Train* sequence at different bit rates.

well as simulation experiments. Two indexes are investigated, algorithm complexity and rate-control precision.

In TM5 algorithm, the final selection of the quantification parameter should consider the MB's spatial activity. The spatial activity of the  $j$ th MB can be calculated by

$$act_j = 1 + \min(vblk_1, vblk_2, \dots, vblk_8) \quad (17)$$

where

$$vblk_n = \frac{1}{64} \times \sum_{k=1}^{64} (P_k^n - P\_mean_n)^2 \quad (18)$$

$$P\_mean_n = \frac{1}{64} \times \sum_{k=1}^{64} P_k^n \quad (19)$$

$P_k$  are the sample values in the  $n$ -th original 8\*8 block.

Based on above analysis, the complexity of TM5 algorithm is  $O(n_1^2)$  according to Formula (18). The complexity of proposed novel algorithm is  $O(n_2^2)$  according to Formulas (8) and (15), where  $n_1$  is the number of MBs in one frame and  $n_2$  is the number of available QPs. Since  $n_1 \gg n_2$ , one can see that the complexity of the improved algorithm reduces greatly compared with the original TM5 algorithm.

To further evaluate the overall performance of the proposed algorithm, simulations are conducted on some standard sequences including *Alex*, *Train* and so on. We mainly investigate the rate-control precision and image quality fluctuation. Here, image quality is computed by *PSNR* and rate-control precision is defined as:

$$\delta = \sqrt{(R - R')^2} / R \quad (20)$$

where  $R$  is the target number of bits of GOP and  $R'$  is the actual generated number of bits.

Figures 4–7 give partial experimental results, where Figure 4 and Figure 5 are the results of rate-control precision of TM5 algorithm and the improved algorithm at different bit rates; Figure 6 and Figure 7 are the results of image quality of the two algorithms. From the experimental results, it can be seen that, compared with the TM5 algorithm, the proposed algorithm can both improve the rate-control precision and reduce the image quality fluctuation.

## 6. Conclusions and Future Work

The paper proposed a novel rate-control algorithm based on the TM5 framework of MPEG-2. The drawback of the target bit allocation method of the original TM5 algorithm is improved and a new rate-distortion model is incorporated. The MB-level rate control is adapted to be frame level. As a result, compared with the original TM5 algorithm, the improved novel algorithm not only can enhance the rate-control precision but also can reduce the

complexity and the fluctuation of decoded image quality. In our future work, the HVS features will be analyzed and is to incorporate into our rate control algorithm to further enhance the subjective visual quality of coded videos.

## 7. Acknowledgments

This work was supported in part by the Natural Science Foundation of Zhejiang Province (No.Y107740); the Open Project Foundation of Ningbo Key Laboratory (No.2007A22002); the Natural Science Foundation of Ningbo (No.2008A610015).

## 8. References

- [1] M. Loren and V. Rahul, "Improved rate control and motion estimation for H.264 encoder," *Proceedings of IEEE International Conference on Image Processing*, pp. 309–312, 2007.
- [2] Z. Z. Chen and K. N. Ngan, "Recent advances in rate control for video coding," *Signal Processing: Image Communication*, pp. 19–38, 2007.
- [3] P. F. Zhao, J. W. Liu, and Q. Li, "A rate control idea and algorithm realization for H.264/AVC," *Computer Engineering*, pp. 233–249, 2006.
- [4] K. Nejat, A. Yucel, and M. M. Russell, "Frame bit allocation for the H.264/AVC video coder via cauchy-density-based rate and distortion models," *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 994–1006, 2005.
- [5] P. Papadimitriou and V. Tsaoussidis, "A rate control scheme for adaptive video streaming over the internet," [Online], <http://www.comp.lancs.ac.uk/~pap-adimp/icc07.pdf>.
- [6] L. Xu, W. Gao, X. Y. Ji, and D. B. Zhao, "Rate control for hierarchical B-picture coding with scaling-factors," [Online], [http://idm.pku.edu.cn/lunwen/2007\\_2008619101121.pdf](http://idm.pku.edu.cn/lunwen/2007_2008619101121.pdf).
- [7] H. Li and Z. Z. Fu, "Improvement research based on H.264 TM5 rate control algorithm," *Journal of Computer Applications*, pp. 749–751, 2007.
- [8] T. Y. Tian, "Research and applications of rate control in video coding," *Journal of University of Electronic Science and Technology of China*, pp. 24–32, 2006.
- [9] W. Yuan, "Research for the ratio control algorithm in H.264," Ph.D thesis, Hefei University of Technology, pp. 9–32, 2006.
- [10] H. J. Bi, *New Generation Video Compression Coding Standard H.264*, Posts and Telecommunications Press, pp. 114–117, 2005.
- [11] Z. J. Zhu, F. Liang, G. Y. Jiang, and M. Yu, "Bit-allocation and rate-control algorithm for stereo video coding," *Journal on Communications*, pp. 15–19, 2007.

# Generation of Multiple Weights in the Opportunistic Beamforming Systems

Guangyue LU<sup>1,2</sup>, Lei ZHANG<sup>2</sup>, Houquan YU<sup>1</sup>, Chao SHAO<sup>2</sup>

<sup>1</sup>*Electronics and Information College, Yangtze University, Jingzhou, China*

<sup>2</sup>*Department of Telecommunications Engineering, Xi'an Institute of Posts and Telecommunications, Xi'an, China*

*E-mail: tonylugy@yahoo.com, chaoshao@xupt.edu.cn*

*Received April 18, 2009; revised April 29, 2009; accepted May 31, 2009*

## Abstract

A new scheme to generate multiple weights used in opportunistic beamforming (OBF) system is proposed to deal with the performance degradation due to the fewer active users in the OBF system. In the proposed scheme, only two mini-slots are employed to create effective channels, while more channel candidates can be obtained via linearly combining the two effective channels obtained during the two mini-slots, thus increasing the multiuser diversity and the system throughputs. The simulation results verify the effectiveness of the proposed scheme.

**Keywords:** Opportunistic Beamforming (OBF), Multiuser Diversity, System Throughputs, Scheduling

## 1. Introduction

With the development of the wireless communication, increasing the spectrum efficiency and data rates is becoming the major task, especially in the downlink case. Multiple-Input-Multiple-Output (MIMO) technique [1] can improve the spectrum efficiency with no need of more bandwidth by employing multiple antennas at both transmitter and receiver. Therefore MIMO technique is becoming one of the most promising techniques in the future communication systems (*e.g.*, LTE, B3G), and coherent beamforming [2] and dirty paper coding [3] are two ways to improving the spectrum efficiency. However the full channel information for all users at the transmitter is required to realize the coherent beamforming and dirty paper coding, which is not realistic with the increasing of the number of the users and antennas because of the waste of the systems resource to feedback the channel information from the receivers to the transmitter.

In wireless communication system, many users are communicating with the base station, and the system throughput can be improved by suitably scheduling (through, *e.g.*, maximum throughput (MAX) scheduling algorithm or proportional fairness (PF) scheduling algorithm) the user with large channel gains to transmit its packets, which is known as the multiuser diversity

(MUD) [4]. In contrast to the channel equalization used in the traditional communication systems to combat the effect of the multipath fading channel on the data transmission, it is the *channel fluctuations* that is the source of the MUD and the MUD will be enlarged with the increase of the dynamic range of the channel fading. The larger the dynamic range of the channel fluctuations, the higher peak of the channels and the larger the multiuser diversity gain. Hence to achieve large MUD requires the large channel dynamic range and the suitable scheduling scheme.

However, the MUD gain will be limited by the small dynamic range of the channel fluctuations due to the availability of light-of-sight (LOS) path and little scattering in the environment and the slowly channel fading compared to the delay constraint of the services. Thus those users with small channel gain and fluctuations may not be scheduled to transmit their packets and their QoS can not be met.

In [5], random fading is induced *purposely* in multiple-input-single-output (MISO) systems when the environment has little scattering and/or the fading is slow to increase the MUD gain of the system by multiplying the transmit data with different weighting factors at each transmitting antenna. When the weighting factors are phase-conjugate with the independent channels from the user to the transmitting antennas, this user is in its

*beamforming configuration state* and its channel peak values occur. When the number of the users in the systems is large enough, the probability that at least one user is in its beamforming configuration state is large and the throughput of the system can approach that of the coherent beamforming with only *partial* channel information (i.e., the overall SNR) feedback. And the scheme in [5] is interpreted as the opportunistic beamforming (OBF).

However, one of the limits of the OBF is the requirement of large number of users in the system simultaneously and the system throughput will be degraded when the number of the users in the system is not too large. When fewer users are active in the system, the MISO system in [5,6] is extended to MIMO in [7], that is, multiple antennas are also employed at the receivers, which equivalently increase the number of virtual active users and, thus, the system throughput. However, the feedback and the costs of each user will be inevitably increased with the increase of the number of the users and the employed receiving antennas. The weighting factors used at the transmitting antennas in [5] are totally random among different time slots. However, since the base station possesses all the users' channels information at current time slot and the previous time slots, the weighting factors can be generated in a pseudo-random manner, that is, the former weighting factors that create beamforming configuration state for one user can be used, in some way, to generate the current former weighting factors only if the coherent time of the channels is large enough [8,9].

Since the random weighting factors strongly affect the channel states, multiple weighting vectors at several mini-slots in one time slot [10] are used to create multiple induced channels, and the one with larger channel gain is selected and the corresponding weighting vector is used as the current weighting vector. The OBF with multiple weighting factors (MW-OBF) can improve the throughput of OBF-CDMA systems. Since several mini-slot are used to 'train' the best weighting factors, some mini-slots and power resources are wasted in MW-OBF.

In [11], two multiple weight OBF schemes tailored for fast fading and slow fading scenarios respectively are investigated and the tight upper bounds of the data rates for both schemes are derived. It is claimed that the faster the fading is, the less weight vectors are desired; and the more users there are, the less weight vectors are desired. To overcome the problem of limited multiuser diversity in a small population, [12] devises a codebook-based OBF (COBF) technique, where the employed unitary matrix changes with time slot to induce larger and faster channel fluctuations in the static channel and to provide further selection diversity to the conventional OBF technique. Compared with [10], the COBF technique reduces

the required number of mini-time slots, and, since it is the size of codebook, not the number of mini-time slots, that determines the amount of supplementary selection diversity, the system throughput can be increased without limitation from the number of mini-time slots. However, the receiver should estimate all of channels from it to the transmitters.

In this paper, a new scheme to generate multiple weights used in OBF is proposed to deal with the performance degradation due to the less number of users in the OBF system. In the proposed scheme, only the equivalent channels at two mini-slots are required to be estimated, as in the normal OBF. The paper is outlined as follow: after the introduction of conventional OBF and MW-OBF in Section 2, the proposed scheme with only two mini-slots to create more channel candidates via linearly combining the two effective channels at the receiver is developed and analyzed in Section 3. Section 4 gives the numerical results to verify the effectiveness of the proposed scheme from different aspects. The paper is concluded in Section 5.

## 2. Conventional OBF and MW-OBF

Assume there are  $N$  transmitting antennas at the base station and one receiving antenna at each user side, the channel gain vector for the  $k$ -th user is  $\mathbf{H}_k(t) = [h_{1k}(t), \dots, h_{Nk}(t)]^T$ , where  $h_{nk}(t)$  ( $n=1, \dots, N$ ) is the channel gain from the  $n$ -th antennas to the  $k$ -th user at time  $t$ . And the transmitting signal  $x(t)$  is multiplied with the weight vector  $\mathbf{V}(t) = \mathbf{e}_\theta^T(t) \boldsymbol{\alpha}(t)$ , where  $\mathbf{V}(t) \in C^{1 \times N}$ , diagonal matrix  $\boldsymbol{\alpha}(t) = \text{diag}(\sqrt{\alpha_1(t)}, \dots, \sqrt{\alpha_N(t)})$  denotes the power allocation on each transmitting antenna, and  $\mathbf{e}_\theta(t) = [e^{j\theta_1(t)}, \dots, e^{j\theta_N(t)}]^T$  is random phase vector applied to the signal,  $\theta_n(t)$  are the independent random variables uniformly distributed over  $[0, 2\pi)$ . In order to preserving the total power,  $\sum_{n=1}^N \alpha_n(t) = 1$ , where random variable  $\alpha_n(t)$  varies from 0 to 1. Then the received signal for the  $k$ -th user is,

$$\begin{aligned} y_k(t) &= \sum_{n=1}^N \sqrt{\alpha_n(t)} e^{j\theta_n(t)} h_{nk}(t) x(t) + z_k(t) \\ &= \mathbf{e}_\theta^T(t) \boldsymbol{\alpha}(t) \mathbf{H}_k(t) x(t) + z_k(t) \\ &\stackrel{\text{def}}{=} \tilde{\mathbf{H}}_k(t) x(t) + z_k(t) \end{aligned} \quad (1)$$

where  $\tilde{\mathbf{H}}_k(t) = \mathbf{e}_\theta^T(t) \boldsymbol{\alpha}(t) \mathbf{H}_k(t) = \mathbf{V}(t) \mathbf{H}_k(t)$  is the equivalent channel (i.e., overall channel) for user  $k$ , and  $z_k(t)$  be the independent and identically distributed AWGN.



From (1), when  $\mathbf{H}_k(t)$  are phase-conjugate with  $\mathbf{e}_\theta(t)$ , that is,  $\theta_n(t) = -\text{angle}(h_{nk}(t))$  ( $n=1, \dots, N$ ),  $\tilde{\mathbf{H}}_k(t)$  are the coherent sum of  $h_{nk}(t)$ , and user  $k$  is in its beamforming configuration state. Thus large channel gain for user  $k$  is obtainable.

In a heavy load system (*i.e.*, the number of active user are large enough), by varying the weights  $\mathbf{V}(t)$ , there is a large possibility that some users are in or nearly in their beamforming configuration states. Using the proportional fair (PF) scheduling algorithm [5], the users with their overall channel SNR near to the peaks are possibly scheduled and the system throughput is approaching to that of the coherent beamforming system.

However, in order to obtain the high throughput by the opportunistic beamforming, a large number of users must exist in each cell. In particular, as the number of transmit antennas of the base station increases, the number of required users grows rapidly. In [9], the conventional OBF is generalized by allowing multiple random weighting vectors at each time slot.

In the multiple weights OBF (MW-OBF) systems, there exist  $Q$  mini-slots in each time slot. During each mini-slots, respectively,  $Q$  known signals multiplied by  $Q$  randomly selected independent weighting vectors  $\mathbf{V}_q(t)$  ( $q=1, \dots, Q$ ) are transmitted. Then, during the  $q$ -th mini-slot, the overall channel gain is

$$\tilde{\mathbf{H}}_{q,k}(t) = \mathbf{V}_q(t) \mathbf{H}_k(t), \quad q=1, \dots, Q \quad (2)$$

Each user measures its overall channel gain,  $|\tilde{\mathbf{H}}_{q,k}(t)|$ , and feeds it back to the base station, then the base station determines the optimum weighting vector,  $w^{opt}(t)$ , for data transmission and the selected user,  $k^*$ ,

$$(k^*, q^{opt}(t)) = \arg \max_{q=1, \dots, Q} \left( \max_{k=1, \dots, K} R_{q,k}(t) \right) \quad (3)$$

$$w^{opt}(t) = w_{q^{opt}(t)}(t) \quad (4)$$

where  $R_{q,k}(t)$  is the transmitted rate for user  $k$  if the  $q$ -th weight vector is used.

### 3. New Scheme to Generate the Multiple Weights

By allowing multiple random weighting vectors at each time slot, the throughput of the MW-OBF scheme is considerably improves compared to the conventional OBF since the employing the weights-selective diversity. However the using of several mini-slots will waste several radio resources and, thus, lower the spectrum efficiency.

In this section, a novel multiple weights generation method is developed by using only two mini-slots at each

time slot. This novel scheme is illustrated with  $N=2$ .

Similar to the MW-OBF, two independent random vectors,  $\mathbf{V}_1(t) = \mathbf{e}_\theta^T(t) \boldsymbol{\alpha}(t)$  and  $\mathbf{V}_2(t) = \mathbf{e}_\phi^T(t) \boldsymbol{\beta}(t)$ , are used at two mini-slots to create two equivalent channels, where  $\boldsymbol{\alpha} = \text{diag}(\alpha_1, \alpha_2)$ ,  $\boldsymbol{\beta} = \text{diag}(\beta_1, \beta_2)$ ,  $\mathbf{e}_\theta^T = (e_{\theta 1}, e_{\theta 2})^T$ ,  $\mathbf{e}_\phi^T = (e_{\phi 1}, e_{\phi 2})^T$ . And the two equivalent channels are, respectively,

$$\tilde{\mathbf{H}}_{eq,k}^{(1)}(t) = \mathbf{V}_1(t) \mathbf{H}_k(t) = \mathbf{e}_\theta^T(t) \boldsymbol{\alpha}(t) \mathbf{H}_k(t)$$

$$\tilde{\mathbf{H}}_{eq,k}^{(2)}(t) = \mathbf{V}_2(t) \mathbf{H}_k(t) = \mathbf{e}_\phi^T(t) \boldsymbol{\beta}(t) \mathbf{H}_k(t)$$

At the receiver, after the estimation of the two equivalent channels, linearly combining them as (the time variable  $t$  is omitted for simplicity in the following),

$$\begin{aligned} \tilde{\mathbf{H}}_{eq,k} &= \tilde{\mathbf{H}}_{eq,k}^{(1)} + b \tilde{\mathbf{H}}_{eq,k}^{(2)} = \mathbf{V}_1 \mathbf{H}_k + b \mathbf{V}_2 \mathbf{H}_k \\ &= \mathbf{e}_\theta^T \boldsymbol{\alpha} \mathbf{H}_k + b \mathbf{e}_\phi^T \boldsymbol{\beta} \mathbf{H}_k = (\mathbf{e}_\theta^T \boldsymbol{\alpha} + b \mathbf{e}_\phi^T \boldsymbol{\beta}) \mathbf{H}_k \end{aligned} \quad (5)$$

where  $\mathbf{H}_k = [h_{1k}, h_{2k}]$ , the complex value  $b$  is the system parameter to be designed as followed.

Denoting

$$\hat{\mathbf{y}}(b) = \mathbf{e}_\theta^T \boldsymbol{\alpha} + b \mathbf{e}_\phi^T \boldsymbol{\beta} \quad (6)$$

then  $\tilde{\mathbf{H}}_{eq,k} = \hat{\mathbf{y}}(b) \mathbf{H}_k$  can be viewed as the OBF channel using weighting factors  $\hat{\mathbf{y}}(b)$ . As in the conventional OBF, to preserve the total transmit power,  $\hat{\mathbf{y}}(b)$  should be normalized as

$$\mathbf{y}(b) = \hat{\mathbf{y}}(b) / |\hat{\mathbf{y}}(b)| \quad (7)$$

Since  $\hat{\mathbf{y}}(b)$  is the function of parameter  $b$ , selecting different  $b$  can resulting in different multiple weighting vectors using only two mini-slots. Then the newly generated channel  $\tilde{\mathbf{H}}_{eq,k}$  is the linear combination of  $\tilde{\mathbf{H}}_{eq,k}^{(1)}$  and  $\tilde{\mathbf{H}}_{eq,k}^{(2)}$ . Suppose that parameter  $b$  is selected from a set with  $W$  elements, then  $W$  new channel can be generated.

In order not to increase the number of multiple operations, suppose  $b$  is selected from the following set,  $\{1, -1, j, -j\}$ , with four elements (*i.e.*,  $W=4$ ). Then six weight vectors can be generated using only two mini-slots, thus improving the spectrum efficiency. Comparing with the original MW-OBF, the proposed scheme needs to estimate the equivalent channels at the two mini-time slots; however, this is easier than the quantized codebook scheme in [11] where channel gains from all users to each antenna must be estimated.

In the proposed scheme, users need feedback its maximum channel gain and the selected parameter  $b$ . Then transmitter schedules the users and calculating the current weights, using (6) and (7) based on the  $b$ ,  $\mathbf{V}_1(t)$

and  $V_2(t)$ .

#### 4. Numerical Results

In this section, we present an extensive set of simulations to verify the effectiveness of the proposed scheme from different aspects. Firstly, since the achievable MUD gain in the system is determined by the dynamic range of the overall channel, which can be described by the probability density function (PDF) of the channels, our simulations depict the PDFs for different schemes. Then, if channels fade very slowly compared to the delay constraint of the application so that transmissions cannot wait until the channel reaches its peak, its QoS cannot be met. Therefore, the channel fluctuation speed, which can be described by the correlation function (CF) of the overall channel, is simulated and given for different schemes. Finally the average throughput of the system for different schemes is simulated for comparison, using both maximum throughput (MAX) scheduling scheme and the PF scheduling scheme.

In the following simulations, we consider two transmit antennas at the base station under the Rician channel with different Rician factor  $\kappa$  and average SNR=0dB. We also suppose the availability of an error-free feedback channel from each user to the base station and the data rate achieved in each time slot is given by the Shannon limit.

##### 4.1. The PDFs and CFs of the Overall Channels

To compare the performance of increasing the dynamic range of the equivalent channels, the PDFs of the chan-

nels are plotted in Figure 1 for Rician channel (with  $\kappa=10$ ) using different schemes, that is, none-OBF, OBF, normal MW-OBF and the proposed scheme. The width of the PDF plot shows the dynamic range of the overall channel. From Figure 1, we can see that the dynamic range of the equivalent channels after OBF and the proposed scheme is much greater than that of the none-OBF, which ensures the larger obtainable MUD gain after OBF and the new MW-OBF scheme. Also comparing the proposed scheme with the normal MW-OBF, OBF and none-OBF, the probabilities that the overall channels have large amplitude are in descending order, which means that the proposed scheme has larger probability to approach high amplitude and, hence, the larger MUD gain.

If the maximum throughput scheduling scheme is employed at the transmitter, the user with the largest channel gain at a time slot will be scheduled to transmit data and the distribution of the peaks of the overall channels will be related to the system throughput directly. Hence, Figure 2 gives the PDFs of the channels' peak for none-OBF, OBF, normal MW-OBF and the proposed scheme, and 10 active users are in the system in the simulation. The four vertical bars, from left to right, indicate the mean values for the four schemes, respectively. The proposed scheme obtains the largest mean values and dynamic range among the four schemes.

Since the fluctuating speed within the time scale of interest is another source of the MUD gain, here we use the normalized correlation function (CF) of the overall channel as the indicator of the fluctuating speed, which is illustrated in Figure 3. And the Rician channels with  $\kappa=30$  are employed in this simulation. For the same

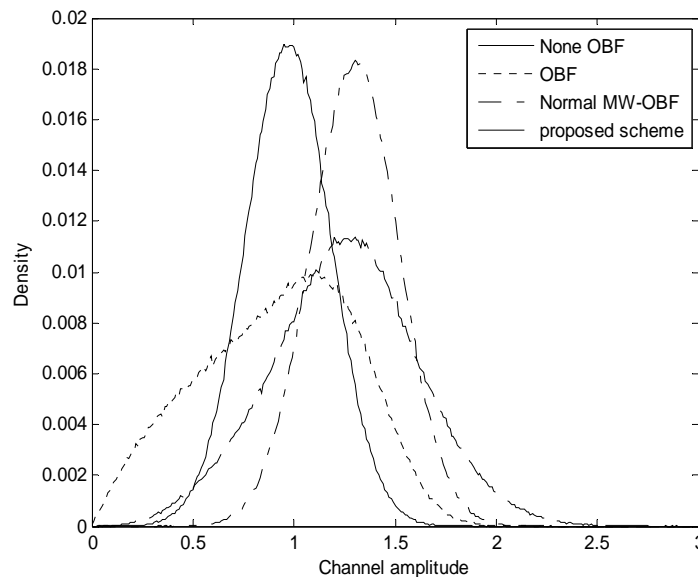


Figure 1. Channels PDFs for Rician channel.

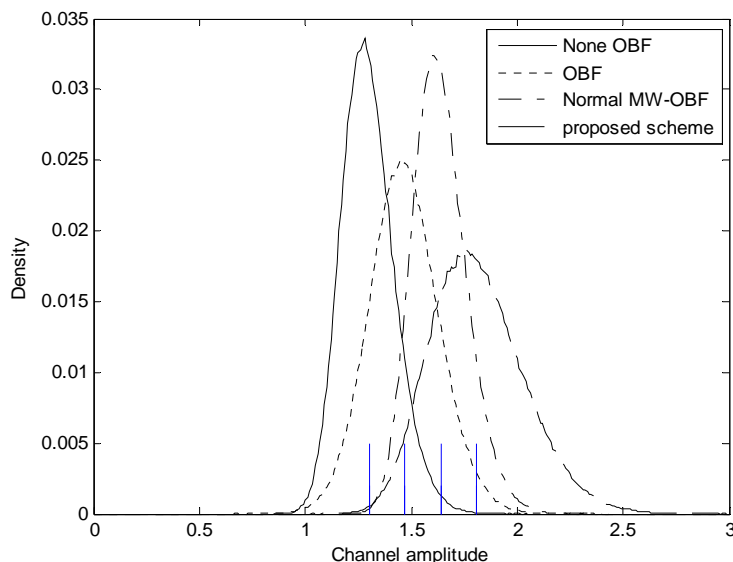


Figure 2. PDFs of channels' peak for Rician channel.

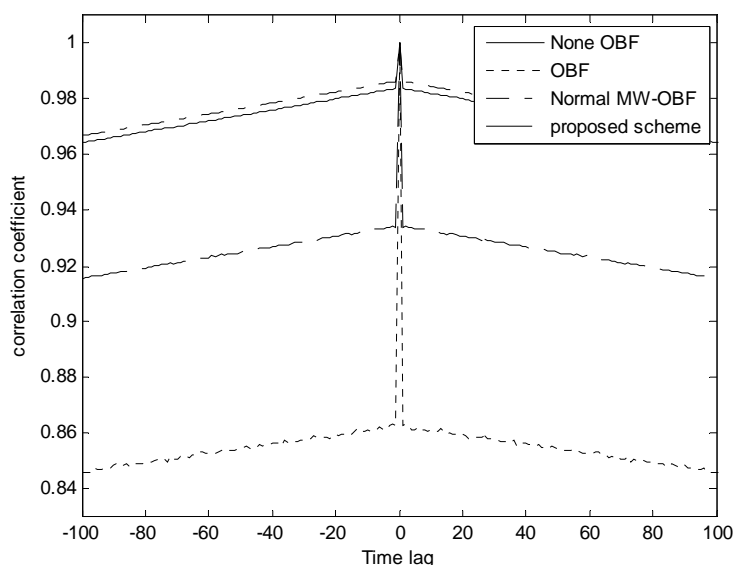


Figure 3. The normalized correlation function of the overall channel.

time lag, the larger the correlation coefficient is, the smaller the fluctuation speed is. So the proposed scheme has less correlation for same time lag, especially for small time lag compared to none-OBF and normal MW-OBF. Since the channels are generated via linearly combining the two equivalent channels, there is correlation among the channels generated in the proposed scheme. So comparing the proposed scheme with OBF, the correlation of the proposed scheme is larger than that of OBF. For example, when the time lag equals  $\pm 1$ , the correlation coefficient of none-OBF, OBF, normal MW-OBF and the proposed scheme are 0.984, 0.98, 0.86 and 0.925, re-

spectively.

From the above simulations, the resulting channels in the proposed scheme have larger dynamic range, larger probability to have high amplitudes, and larger fluctuating rate. We, therefore, can expect that the proposed scheme can obtain larger MUD gain, which will be illustrated in the following simulations.

#### 4.2. The System Average Throughput for Different Schemes

The simulating parameters are same as those in [10]. The

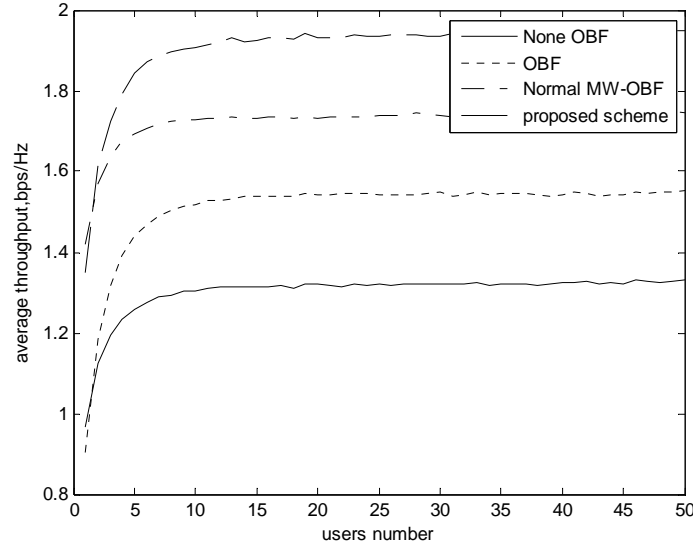


Figure 4. Average throughput using the PF scheme.

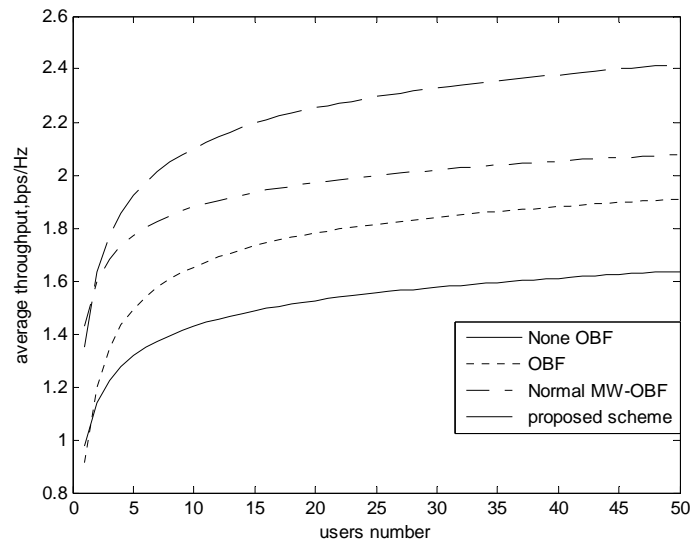


Figure 5. Average throughput using MAX scheduling scheme.

Rician channel with  $\kappa = 10$ . Six mini-slots are used to generate six equivalent channels in MW-OBF, whereas two mini-slots are used in the proposed scheme to create two overall channels, and four additional channels are generated via linearly combining the available two overall channels.

Figures 4 and 5 illustrate the average throughput of the system using PF and MAX scheme for different schemes, respectively. The results show that, in both scheduling schemes, the average throughput are improved greatly, especially when the system with small number of users in MW-OBF and the proposed scheme. Meantime, the proposed scheme has larger throughput than MW-OBF.

#### 4.3. Throughput Variation with the Rician Factors

Finally we study the performance variation of different schemes with the Rician factor  $\kappa$ , that is, to investigate the influence of the light of sight (LOS) on the system throughput. From the Figure 6, it can be seen that with the increase of  $\kappa$ , the throughput for all scheme degrades because the throughput rely on the peak values of the instant overall channel. When the  $\kappa$  factor increases, the channel fluctuations are reduced and the peak values of the instant overall channel are reduced, too. Compared with normal OBF, the proposed scheme can be improved the throughput, for example, for  $\kappa = 10$ ,

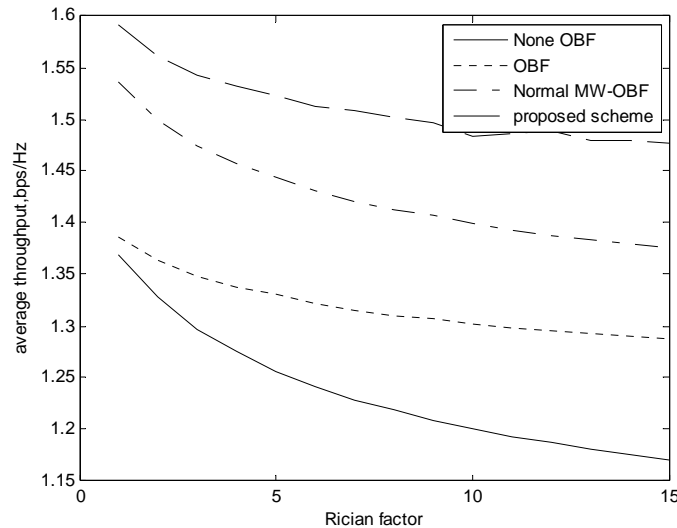


Figure 6. Throughput versus Rician factor of the channel.

more than 10% throughput enhancement can be obtained.

## 5. Conclusions

A new simple scheme to generate multiple weights used in opportunistic beamforming (OBF) system is proposed in this paper to deal with the performance degradation due to the fewer users in the OBF system. Only two mini-slots are employed to create effective channels, while more channel candidates can be obtained via linearly combining the two effective channels at the receiver side, thus increasing the multiuser diversity and the system throughputs. The simulation results show that the throughput can be improved using the proposed scheme.

## 6. Acknowledgements

This work is supported by Program for New Century Excellent Talents in University (NCET-08-0891), the Natural Science Foundation of China under the grant No. 60602053, and the Natural Science Foundation of Shaanxi Province under the grant No. 2007F02.

## 7. References

- [1] E. Telatar, "Capacity of multi-antenna Gaussian channels," *European Transactions on Telecommunications*, Vol. 10, No. 6, pp. 585–596, 1999.
- [2] F. Rashid-Farrokhi, K. J. R. Liu, and L. Tassiulas, "Transmit beamforming and power control for cellular wireless systems," *IEEE Journal in Selected Areas on Communications*, Vol. 16, No. 8, pp. 1437–1450, August 1998.
- [3] M. Costa, "Writing on dirty paper," *IEEE Transactions on Information Theory*, Vol. 29, No. 3, pp. 439–441, May 1983.
- [4] R. Knopp and P. A. Humblet, "Information capacity and power control in single cell multiuser communications," In *Proceedings of IEEE International Conference on Communications*, pp. 331–335, 1995.
- [5] P. Viswanath, D. N. C. Tse, and R. Laroia, "Opportunistic beamforming using dumb antennas," *IEEE Transactions on Information Theory*, Vol. 48, No. 6 pp. 1277–1294, June 2002.
- [6] M. Sharif and B. Hassibi, "On the capacity of MIMO broadcast channels with partial side information," *IEEE Transactions on Information Theory*, Vol. 51, No. 2, pp. 506–522, February 2005.
- [7] W. Zhang and K. B. Letaief, "MIMO broadcast scheduling with limited feedback," *IEEE Journal in Selected Areas on Communications*, Vol. 25, No. 7, pp. 1457–1467, July 2007.
- [8] M. Kountouris and D. Gesbert, "Memory-based opportunistic multi-user beamforming," In *Proceedings of International Symposium on Information Theory*, pp. 1426–1430, September 2005.
- [9] I. R. Baran and B. F. Uchoa-Filho, "Enhanced opportunistic beamforming for Jakes-correlated fading channels," In *Proceedings of International Telecommunications Symposium*, pp. 1024–1029, Fortaleza, Ceara, September 2006.
- [10] H.-M. Kim, S. C. Hong, and S. S. Ghassemzadeh, "Opportunistic beamforming based on multiple weighting vectors," *IEEE Transactions on Wireless Communications*, Vol. 4, No. 6, pp. 2683–2687, November 2005.
- [11] M. Zeng, J. Wang, and S. Q. Li, "Rate upper bound and optimal number of weight vectors for opportunistic beamforming," In *proceedings of IEEE Vehicular Technology Conference*, Fall, pp. 661–665, September 30 2007–October 3, 2007.
- [12] J. Kang, I. K. Choi, D. S. Kwon, and C. Y. Lee, "An opportunistic beamforming technique using a quantized codebook," In *proceedings of IEEE Vehicular Technology Conference*, pp. 1647–1651, Spring, 2007.

# The Effect of Notch Filter on RFI Suppression

Wenge CHANG, Jianyang LI, Xiangyang LI

*School of Electronics Science and Engineering, National University of Defense Technology, Changsha, China*

*E-mail: changwenge@sina.com*

*Received May 22, 2009; revised May 31, 2009; accepted June 10, 2009*

## Abstract

Radio Frequency Interference (RFI) suppression is an important technique in the ultra-wideband synthetic aperture radar (UWB SAR). In this paper, we mainly analyze the performance of a notch filter for RFI suppression. The theoretical output from notch filter is presented based on RFI signal's narrowband property. The research conclusion shows that the notch filter has significant effect on sidelobes of the system response, which might be considered to be false targets, however it has little effect on the resolution of the system response. The theoretical result is verified by simulation and experimental data processing both in one dimension (range dimension) and in two dimensions (range and azimuth dimension).

**Keywords:** RFI Suppression, Matched Filter, Notch Filter, SAR

## 1. Introduction

The dual requirement of a low radar frequency for foliage and/or ground penetration and a wide radar bandwidth for high resolution in wideband radar systems leads to radar operating in frequency bands occupied by other radio systems, such as TV and radio communications. As a result, Radio Frequency Interference (RFI) appears in the received radar signal [1–3]. In ultra-wideband synthetic aperture radar (UWB SAR), the RFI energy is spread over the whole image scene, displaying artefacts and masking targets, especially in low SNR areas [1–5]. Any subsequent processing in UWB SAR (such as target classification, interferometry, etc.) would be degraded by the presence of RFI. Thus, RFI suppression is an important technique in UWB SAR signal processing.

The common approach to RFI suppression is to examine the spectrum of the contaminated signal, identify the interference spikes which usually have greater power than the radar echo signal, and subsequently remove these spikes with the help of a notch filter [1]. The notch concept is effective for RFI suppression. Having studied the notch filter based on least-mean-squared estimation and tested with real RFI data, T. Koutsoudis and L. Lovas suggested that the notch filter can produce an adverse impact on the SAR performance (such as reducing image intensity, range resolution and creating loss in the target's signal to noise ratio) but no theoretic results were presented in reference to this claim [2].

In this study, the performance of the notch filter is theoretically analysed. Firstly, the transmitted radar signal model, the received radar signal model and the RF interference model are presented. A matched filter with notches is designed. Secondly, the contaminated signal (the sum of radar echo and RFI) is fed to the matched filter, and the theoretical output of the filter is derived. The theoretical result shows that the output of the matched filter with notches is influenced by the notches' width and carrier and, furthermore that the notch filter has little impact on range resolution but significant impact on sidelobes. The simulation is studied both in the case of one RF interference existing as well as multiple RF interference existing. Finally, the matched filter with notches is applied to the experimental data acquired by an airborne UWB SAR, and the validity of the theoretical result is tested.

This paper is arranged as follows: In part 2, the models are built and the output from the notch filter is theoretical deduced. In part 3, the simulation in range is carried out and the result is shown to be consistent with theoretical result. In part 4, the SAR image simulation as well as the experimental UWB SAR image is processed in order to test and verify the validity of the theoretical result. Finally, the results are summarised in the conclusion.

## 2. Modelling and Theoretical Deducing

Two assumptions are made in order to make analysis

more convenient. Firstly, it is assumed that the transmitted and received signal of UWB SAR is a base-band linear frequency modulated pulse as follows:

$$s(t) = \text{rect}\left(\frac{t}{T}\right)e^{j\pi k^2}, k = B/T; \quad (1)$$

$$S(\omega) = \text{rect}\left(\frac{\omega}{2\pi B}\right)e^{-j\frac{\omega^2}{4\pi k}} \quad (2)$$

where  $B$  is the bandwidth,  $T$  is the width of the transmitted pulse,  $S(\omega)$  is the spectrum of  $s(t)$ .

Secondly, it is assumed that only one narrow band RFI signal exists with carrier  $\omega_1$  and bandwidth  $b_1$ .

The matched filter with a notch at  $\omega_1$  is designed for suppressing FRI with carrier  $\omega_1$  and bandwidth  $b_1$ . The matched filter in spectrum is shown in Equation (3).

$$H(\omega) = (1 - \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right))S^*(\omega) = S^*(\omega) - \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{j\frac{\omega^2}{4\pi k}} \quad (3)$$

The received radar signal includes the echo signal (the reflected transmitted signal from the target), the RFI signal and thermal noise. Assuming only one RFI  $s_{RFI}(t)$  exists and that thermal noise can be ignored, the received signal  $s_r(t)$ , and its spectrum  $S_r(\omega)$ , can be written as

$$s_r(t) = a \cdot s(t - \tau) + s_{RFI}(t) \quad (4)$$

$$S_r(\omega) = aS(\omega)e^{-j\omega\tau} + S_{RFI}(\omega) \quad (5)$$

where  $a$  is a constant coefficient. The received signal is processed with the matched filter. The output in spectrum is given by Equation (6).

$$\begin{aligned} Y(\omega) &= S_r(\omega)H(\omega) = (aS(\omega)e^{-j\omega\tau} \\ &+ S_{RFI}(\omega)) \cdot (S^*(\omega) - \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{j\frac{\omega^2}{4\pi k}}) \\ &= a \cdot \text{rect}\left(\frac{\omega}{2\pi B}\right)e^{-j\omega\tau} + S^*(\omega)S_{RFI}(\omega) \\ &\quad - aS(\omega)e^{-j\omega\tau} \cdot \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{j\frac{\omega^2}{4\pi k}} \\ &\quad - a \cdot \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{j\frac{\omega^2}{4\pi k}} \cdot S_{RFI}(\omega) \end{aligned} \quad (6)$$

where:

$$\begin{aligned} &aS(\omega)e^{-j\omega\tau} \cdot \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{j\frac{\omega^2}{4\pi k}} \\ &= a \cdot \text{rect}\left(\frac{\omega}{2\pi B}\right)e^{-j\frac{\omega^2}{4\pi k}} \cdot e^{-j\omega\tau} \cdot \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{j\frac{\omega^2}{4\pi k}} \\ &= a \cdot \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right) \cdot e^{-j\omega\tau} \end{aligned} \quad (7)$$

$$S^*(\omega)S_{RFI}(\omega) = \text{rect}\left(\frac{\omega}{2\pi B}\right)e^{j\frac{\omega^2}{4\pi k}} \cdot S_{RFI}(\omega) \quad (8)$$

Thus:

$$\begin{aligned} Y(\omega) &= a \cdot \text{rect}\left(\frac{\omega}{2\pi B}\right)e^{-j\omega\tau} - a \cdot \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{-j\omega\tau} \\ &\quad + \text{rect}\left(\frac{\omega}{2\pi B}\right)e^{j\frac{\omega^2}{4\pi k}} \cdot S_{RFI}(\omega) - \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{j\frac{\omega^2}{4\pi k}} \cdot S_{RFI}(\omega) \end{aligned} \quad (9)$$

Based on the assumption that the RFI signal  $s_{RFI}(t)$  is a narrow band signal with carrier  $\omega_1$  and bandwidth  $b_1$ , we have

$$\text{rect}\left(\frac{\omega}{2\pi B}\right)e^{j\frac{\omega^2}{4\pi k}} \cdot S_{RFI}(\omega) \approx \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{j\frac{\omega^2}{4\pi k}} \cdot S_{RFI}(\omega) \quad (10)$$

Considering Equations (10) and (9), the output of the matched filter with a notch is given by

$$Y(\omega) = a \cdot \text{rect}\left(\frac{\omega}{2\pi B}\right)e^{-j\omega\tau} - a \cdot \text{rect}\left(\frac{\omega - \omega_1}{2\pi b_1}\right)e^{-j\omega\tau} \quad (11)$$

The output in the time domain of the matched filter with a notch is given by

$$y(t) = aB \cdot \text{Sa}(B(t - \tau)) - ab_1 \cdot \text{Sa}(b_1(t - \tau))e^{-j\omega_1(t + \tau)} \quad (12)$$

From Equation (12) we find that the output of the matched filter with a notch has two parts: one is the desired part and the other one is the undesired part. Here we call the undesired part *clutter*, which is a high frequency oscillating signal modulated by a wide pulsed *sinc*-function, determined by parameters  $\omega_1$ ,  $b_1$ . This clutter influences the peak-sidelobe-ratio (PSLR), the integrated-sidelobe-ratio (ISLR) of radar system response.

Here the signal-to-clutter-ratio (SCR) is utilised to describe the influence as follows:

$$SCR = 20\log(b_1 / B) \quad (13)$$

It is easy to expand Equation (12) for multiple RF interference signals. Assuming there are  $n$  RFI signals with carriers and bandwidths being  $\omega_n, b_n$  ( $n=1,2,\dots,n$ ) respectively.

The matched filter with  $n$  notches is given by:

$$\begin{aligned} H(\omega) &= (1 - \sum_{k=1}^n \text{rect}\left(\frac{\omega - \omega_k}{2\pi b_k}\right))S^*(\omega) \\ &= S^*(\omega) - \sum_{k=1}^n \text{rect}\left(\frac{\omega - \omega_k}{2\pi b_k}\right) \cdot e^{j\frac{\omega^2}{4\pi k}} \end{aligned} \quad (14)$$

The output spectrum of the matched filter is:

$$Y(\omega) = a \cdot \text{rect}\left(\frac{\omega}{2\pi B}\right) e^{-j\omega\tau} - a \cdot \sum_{k=1}^n \text{rect}\left(\frac{\omega - \omega_k}{2\pi b_k}\right) e^{-j\omega\tau} \quad (15)$$

Corresponding time domain expression is:

$$y(t) = aB \cdot \text{Sa}(B(t-\tau)) - a \cdot \sum_{k=1}^n b_k \cdot \text{Sa}(b_k(t-\tau)) e^{-j\omega_k(t+\tau)} \quad (16)$$

For  $n$  RFI signals, the output of the matched filter with  $n$  notches also consists of two parts: one is desired, the other one is undesired (which is the sum of  $n$  high frequency oscillating signal modulated by wide pulsed *sinc*-function, determined by parameters  $\omega_n$ ,  $b_n$   $n=1, 2, \dots, n$ ).

SCR is no longer suitable for describing the influence of the notch filter under these conditions. For the clutter, being coherent, might accumulate, some of the sidelobes would be higher than the case when only one RFI signal exists.

From the Equation (12) and Equation (16), we have the conclusion that the notch filter affects the output of matched filter in following ways:

- It produces the undesired output that is an oscillating signal modulated by a *sinc*-function. The undesired output influences PSRL and ISLR of radar system response.
- It has little influence on the resolution of radar system response.

### 3. Simulation

We have two simulation steps in order to demonstrate the validity of Equation (16) with only one RFI and multiple RFI signals existing.

#### 3.1. Only One RFI Existing

It is assumed that the base-band LFM signal has 20us pulse-width with 200MHz bandwidth, and that the RFI carrier is 30MHz and the bandwidth is 8MHz. Under the above conditions, the spectrum of the matched filter with a notch is shown in Figure 1. The output of the matched filter with a notch is shown in Figure 2. In the figure the ideal output, the output of the ideal matched filter for a ideal LFM input, is as dotted line and the output of the matched filter with RFI suppression for the mixed signal is as solid line. Comparing two outputs we find that regular spikes appear in the output with RFI suppression.

To describe the effect of the notch filter in a clear way, the output is separated two parts, as the Equation (12), described as in Figure 3. The solid line is the desired output and the dotted line is undesired part. From the figure, the SCR is about  $-27.6\text{dB}$ , being consistent with the theoretical result of  $-27.96\text{dB}$  derived from Equation (13).

Figure 4 is the output of ideal matched filter for mixed

signal. From the figure, the sidelobes rise, esp., the sidelobes being far from the mainlobe, as a result the clutter rise greatly.

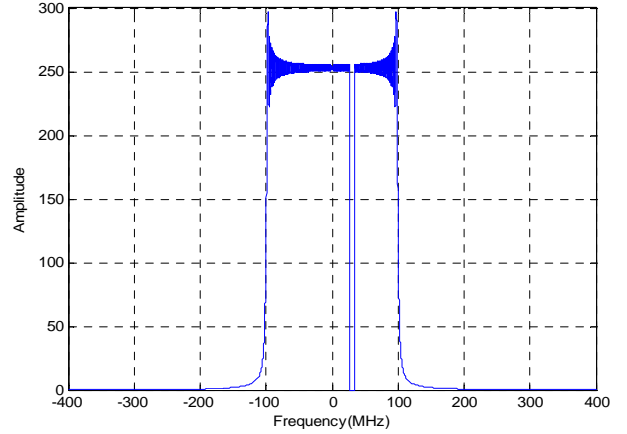


Figure 1. Matched filter with a notch.

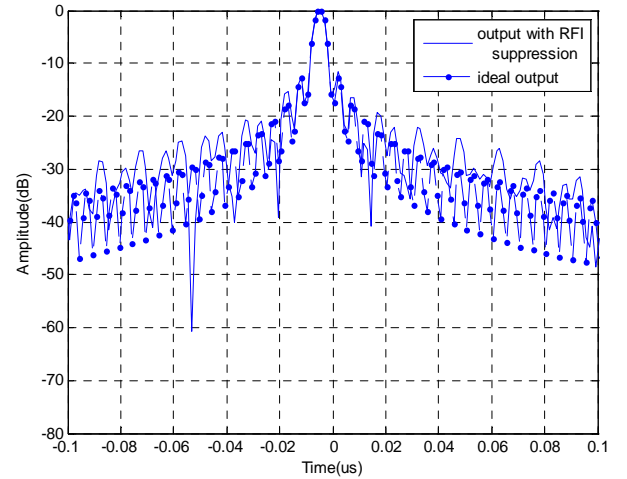


Figure 2. The output of matched filter.

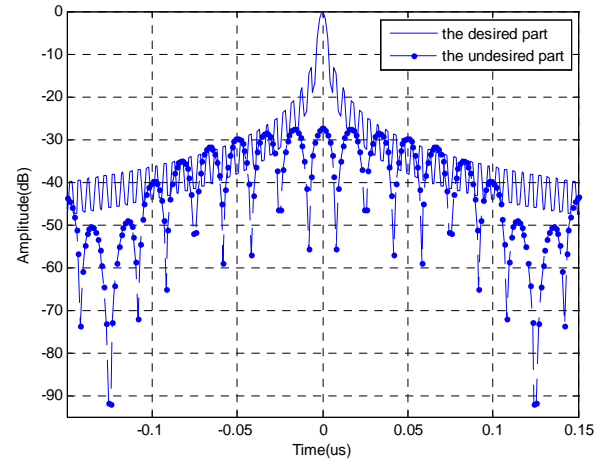


Figure 3. The separated form of output.



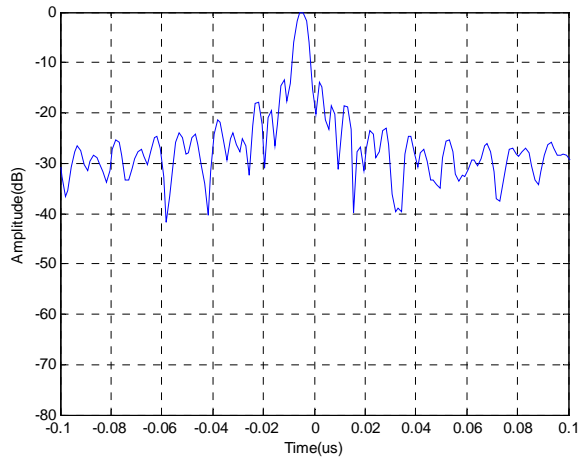


Figure 4. The output for mixed signal.

From the figures we find that the notch filter do have significant effect on RFI suppression, however, it has side effect that the sidelobes raise greatly.

### 3.2. Multiple RFI Existing

It is assumed that the base-band LFM signal is as same as above. Meanwhile, three existing RFI signals are assumed and the RFI signals' carriers and bandwidths are (−50MHz, 8MHz), (30MHz, 8MHz), and (50MHz, 8MHz) respectively. Under these conditions, the spectrum of the matched filter with notches is shown in Figure 5.

The output of the matched filter with notches is shown in Figure 6. In the figure the ideal output is as dotted line and the output of the matched filter with RFI suppression for the mixed signal is as solid line. Comparing two outputs we also find that regular spikes appear in the output with RFI suppression.

The output is also separated two parts, as the Equation

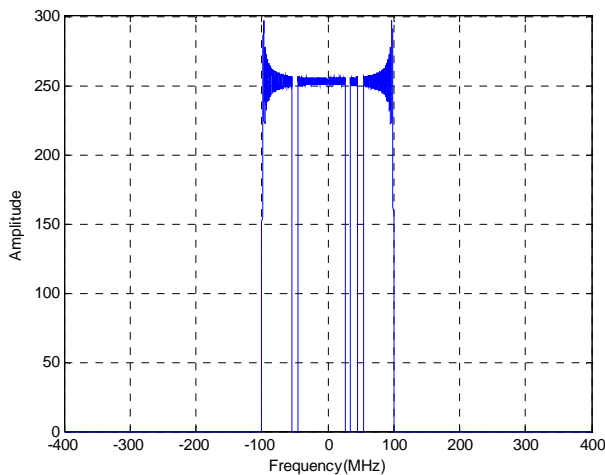


Figure 5. Matched filter with notches.

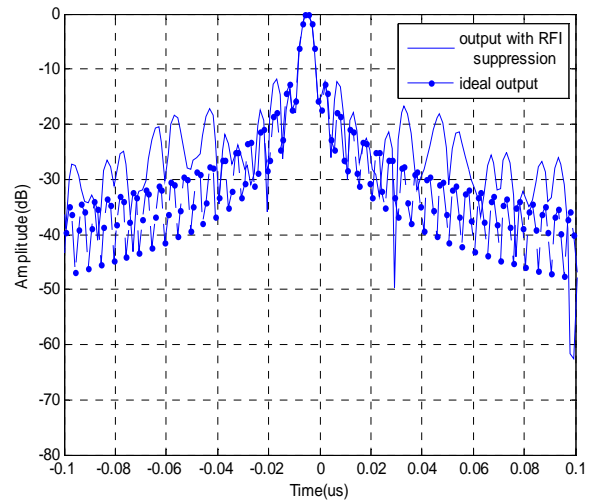


Figure 6. The output of matched filter with notches.

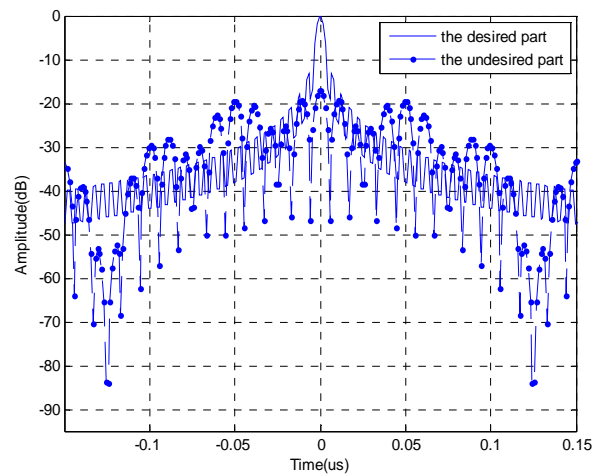


Figure 7. The separated form of output.

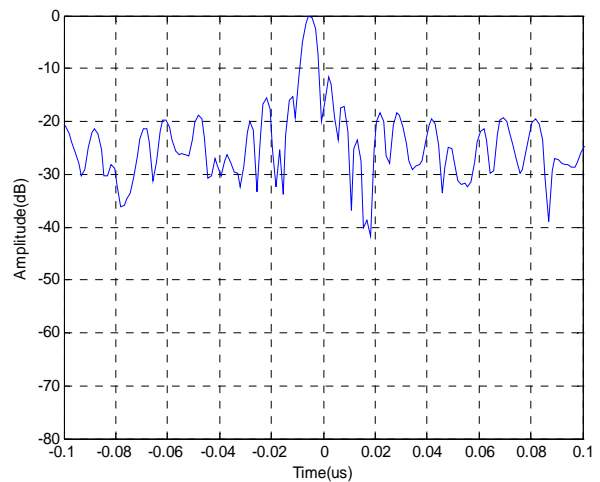


Figure 8. The output for mixed signal.

(16), described as in Figure 7. It can be observed that the *SCR* is significantly higher than in the case of only one existing RFI. We measure the highest *SCR* as  $-18.2\text{dB}$ , which is consistent with the theoretical result  $20\text{Lg}(3b/B) = -18.42\text{dB}$  derived from Equation (13).

Figure 8 is the output of ideal matched filter for mixed signal. From the figure, we have the same conclusion as the former: sidelobes rise, esp., the sidelobes being far from the mainlobe, causing the clutter rise greatly.

We also find that the notch filter do have significant effect on RFI suppression, however, it has side effect. In summary, the notch filter has significant effect on RFI suppression, but it can produce an adverse impact that the sidelobes rise on the output of the matched filter. Meanwhile we find that notch filter has little impact on range resolution.

#### 4. Real Data Processing

To check the validity of the theoretical conclusion we apply the matched filter with notch(s) to the real data

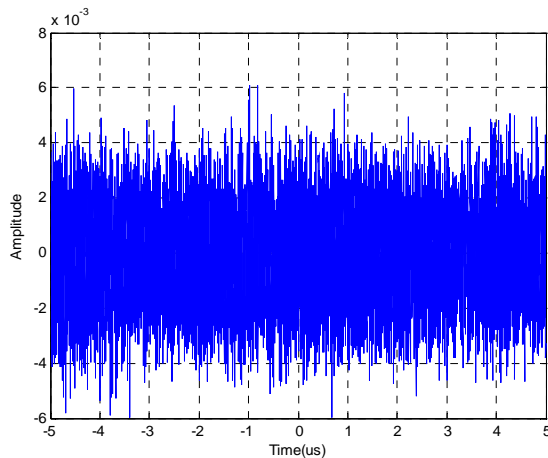


Figure 9. The received signal.

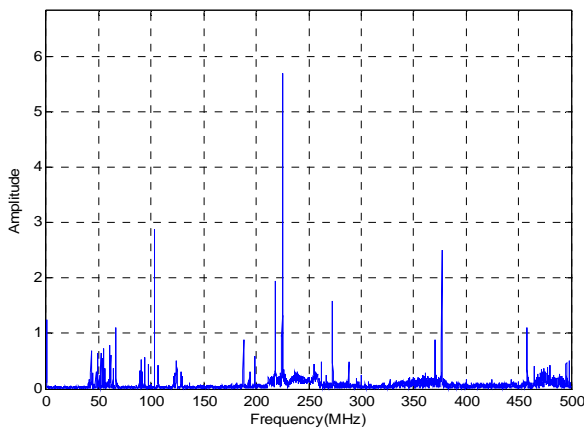


Figure 10. The spectrum of received signal.

from a specified system. In the system, the transmitted signal, which is LFM signal with  $10\mu\text{s}$  pulse-width,  $50\text{MHz}$  bandwidth and  $235\text{MHz}$  centre frequency, is generated by AWG520. This signal is fed to an antenna which is Archimedes screw form to radiation. Meanwhile a same form antenna is used to receive. The distance between two antennas is  $8\text{m}$ . The received signal is amplified, and then is sampled by TDS784D Digital Oscilloscope. The received signal is as Figure 9. The spectrum is as Figure 10. From the Figures, we can see there are two RFI signals in the bandwidth of  $[210\text{MHz} \sim 260\text{MHz}]$ , but only one RFI is relatively strong. So we use the notch filter to suppress this stronger one. The notch filter's carrier and bandwidth are  $221.5\text{MHz}$ ,  $6.5\text{MHz}$  respectively.

The output of matched filter without notch is shown in Figure 11. It is clear that RFI has intense influence that sidelobes, whatever are far or near from mainlobe, are raised to about  $-13\text{dB}$  on the system response.

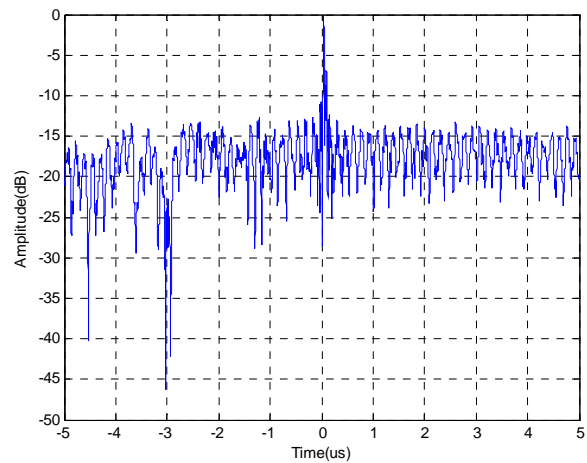


Figure 11. The output of matched filter without notch.

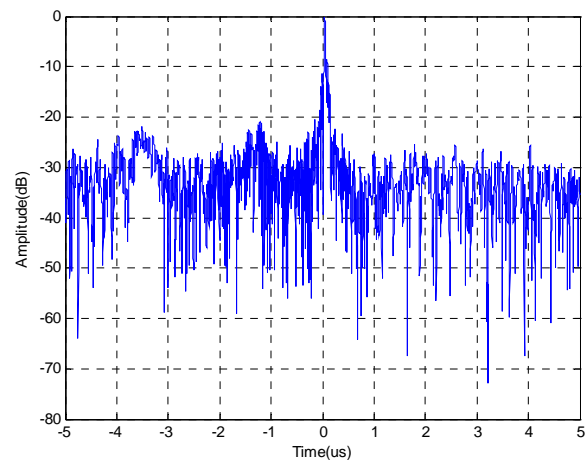


Figure 12. The output of matched filter with notch.

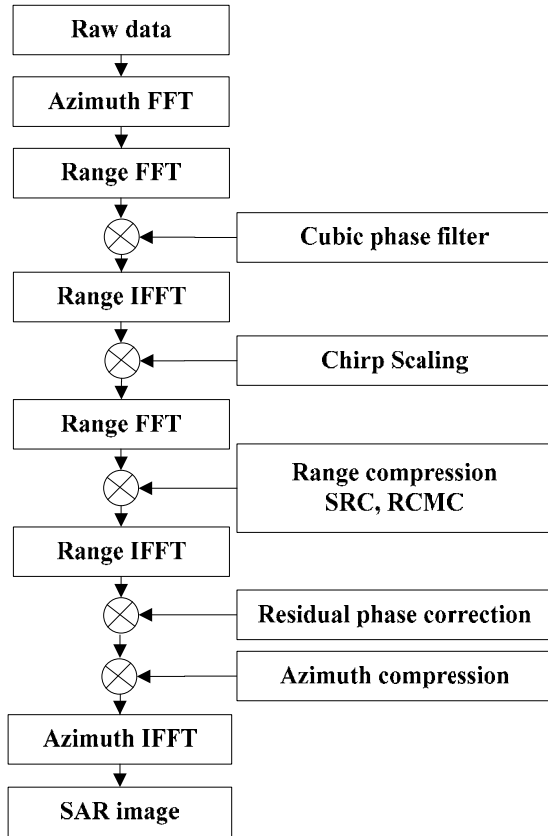


Figure 13. The diagram of NCS image algorithm.

The output of matched filter with notch is shown in Figure 12. It can be seen that the influence of RFI have been decreased greatly, and the sidelobes level is less than  $-21\text{dB}$ . However, the regular spikes caused by the notch filter appear.

## 5. Imaging

The statement that the notch filter can suppress RFI but it has negative impact on the matched filter's output has already been demonstrated in the range signal processing. Furthermore, we will demonstrate the theoretical conclusion both in range and azimuth signal processing, that is, in SAR image processing.

Through illuminating the ground with coherent radiation and measuring the echo signals, SAR can produce two dimensional imageries with high resolution of the ground surface. Range resolution is accomplished through range gating. Fine range resolution can be accomplished by using pulse compression techniques. The azimuth resolution depends on antenna size and radar wavelength. Fine azimuth resolution is enhanced by taking advantage of the radar motion in order to synthesize a larger antenna aperture.

The statement for the notch filter being used in SAR processing is also that significantly higher sidelobes would blur the SAR image where RFI signal exist.

SAR image for simulation and experimental UWB SAR data is used to further demonstrate the influence of the notch filter. The procedure of image processing (whatever the simulation or experimental data processing is) is as follows:

- 1) generating or receiving the radar echo signal mixed with the RFI signal;
- 2) analysing the echo signal spectrum and identifying the RFI signal, including the RFI signal's carriers and bandwidth;
- 3) designing the matched filter with notches, as in Equation (14);
- 4) imaging with the NCS image algorithm [6]. The procedure is illustrated in Figure 13.

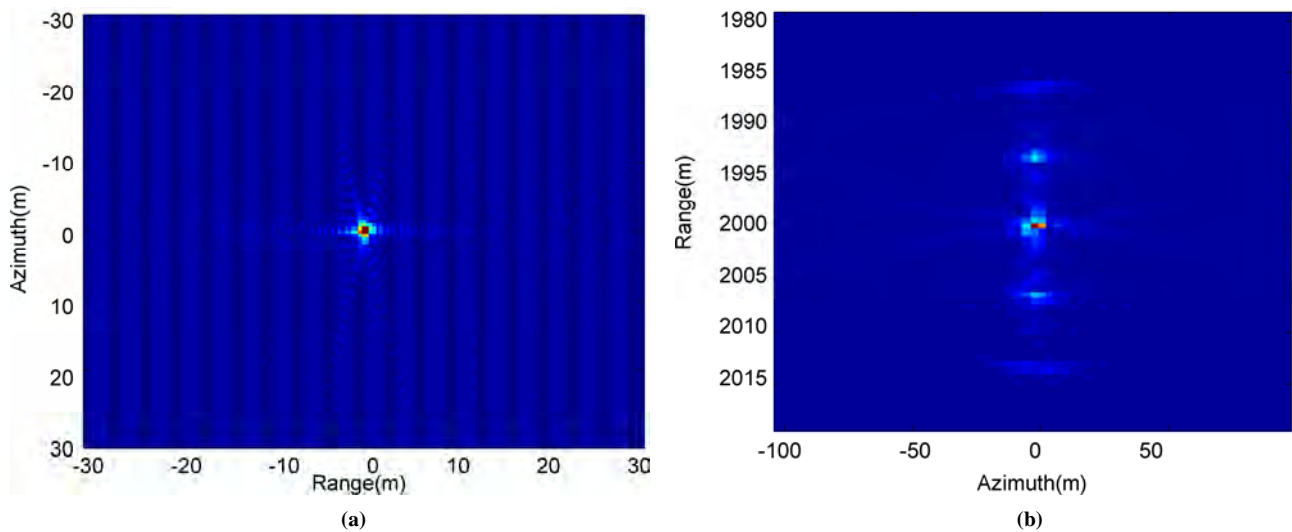


Figure 14. Imaging for four RFIs with 4MHz bandwidth (a: Without RFI suppression, b: With RFI suppression).

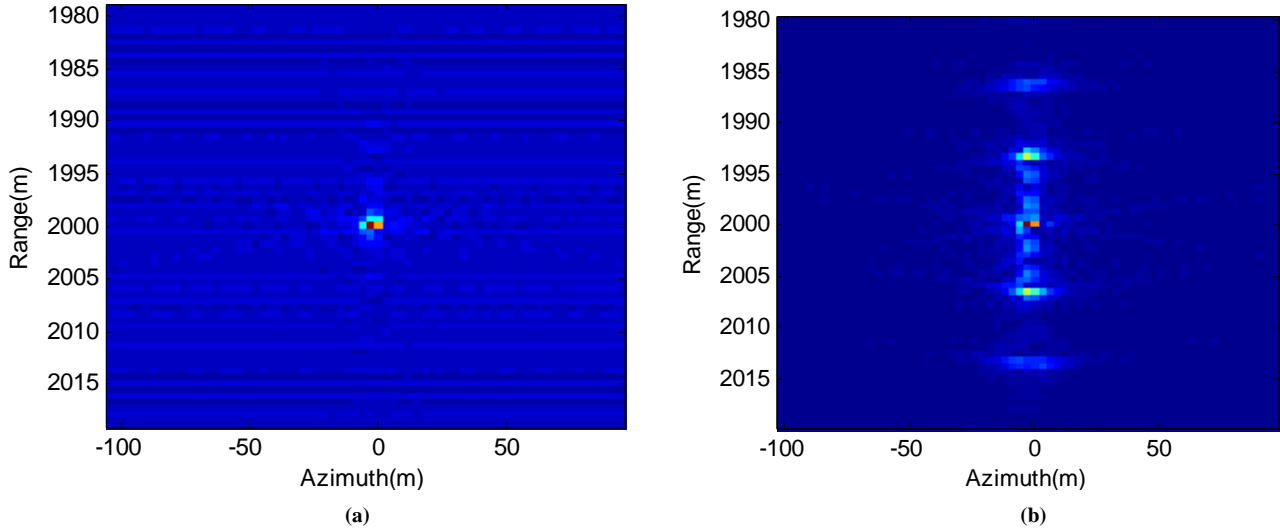


Figure 15. Imaging for four RFIs with 8MHz bandwidth (a: Without RFI suppression, b: With RFI suppression).

### 5.1. Imaging for Simulation

In order to study the effect of the notch filter, we have LFM signal with 20 $\mu$ s pulse-width and 200MHz bandwidth to mix with four RFI signals to form the radar echo.

Firstly, the bandwidths of all the RFI signals are assumed to be the same as 4MHz, and the carriers are 22MHz, 44MHz, 66MHz and 88MHz respectively. After the filter with the four notches and the SAR image is processed, the result is shown in Figure 14. Figure 14(a) is the figure without any RFI suppression, so the regular texture caused by RFI can be seen clearly and Figure 14(b) is the figure with RFI suppression using the notch filter, where the regular texture disappears, but the false targets in range, caused by the notch filter, appear.

Then all the RFI signal bandwidths are assumed to be 8MHz whereas the other conditions are kept as above. The imaging result is shown in Figure 15. Figure 15(a) is the figure without any RFI suppression and Figure 15(b) is the figure with RFI suppression using the notch filter.

Comparing Figure 14 to Figure 15, we find that RFI certainly blur the SAR image without RFI suppression, whereas the false targets appear in the SAR image with RFI suppression using the notch filter. From the simulation we find that the wider bandwidth the RFI signals have, the stronger the false targets appear in the image. This observation is consistent with the theoretical conclusion.

The influence of the RFI quantities has been studied, and with Figure 14 and Figure 15, the numerical values are measured and described in Table 1 and Table 2, where Table 1 is the result without any RFI suppression and Table 2 is the result with RFI suppression using the

notch filter. In the table, the term ( $N \times$ RFI with BMHz) means that there are  $N$  RFI signals with BMHz bandwidth in the radar echo.

We draw further conclusions from Table 1 and Table 2, namely that the notch filter has significant effect on sidelobes, as the quantities and the RFI bandwidth increase, the sidelobes rise too. Raised sidelobes result in the SAR image having more false targets.

We also have the conclusion from Table 1 and Table 2 that the notch filter has little effect on the SAR resolution.

### 5.2. Imaging for Experimental UWB SAR Data

The matched filter with notches is applied to experimental UWB SAR data to verify the simulation and theoretical

Table 1. SAR performance without RFI suppression.

	Resolution(m)	PSLR(dB)	ISLR(dB)
2 $\times$ RFI with 4MHz	0.66(r) $\times$ 0.9(a)	-11.32	-9.71
2 $\times$ RFI with 8MHz	0.66(r) $\times$ 0.9(a)	-11.8	-9.71
8 $\times$ RFI with 4MHz	0.66(r) $\times$ 0.9(a)	-11.23	-9.15
8 $\times$ RFI with 8MHz	0.66(r) $\times$ 0.9(a)	-11.59	-9.74

\*Note: ISLR is calculated only in range dimension.

Table 2. SAR performance with RFI suppression.

	Resolution(m)	PSLR(dB)	ISLR(dB)
2 $\times$ RFI with 4MHz	0.66(r) $\times$ 0.9(a)	-11.12	-6.04
2 $\times$ RFI with 8MHz	0.66(r) $\times$ 0.9(a)	-7.84	-2.15
8 $\times$ RFI with 4MHz	0.66(r) $\times$ 0.9(a)	-11.12	-4.12
8 $\times$ RFI with 8MHz	0.66(r) $\times$ 0.9(a)	-9.06	1.1

\*Note: ISLR is calculated only in range dimension.

conclusion. The UWB SAR works at UHF frequency band with a bandwidth of 200MHz. Furthermore, The UWB SAR is a subsystem of the multi-frequency-band SAR (MFB SAR) developed in 2005 by the National University of Defense Technology (NUDT) of China, associated with the East China Research Institute of Electronic Engineering (ECRIEE). The MFB SAR system installed on an Y7 aeroplane is capable of operating simultaneously in four frequency bands. Its flight test was performed in January 2005 in an area near the Sanya, Hainan, China. The sky is full of maritime radio and TV signals. Figure 16 is the averaged range signal spectrum of the received base-band signal. This spectrum is real, so we can see that there are at least 6 radio frequency interference signals in the band and their power vary from 10dB to 20dB greater than the radar echo signal.

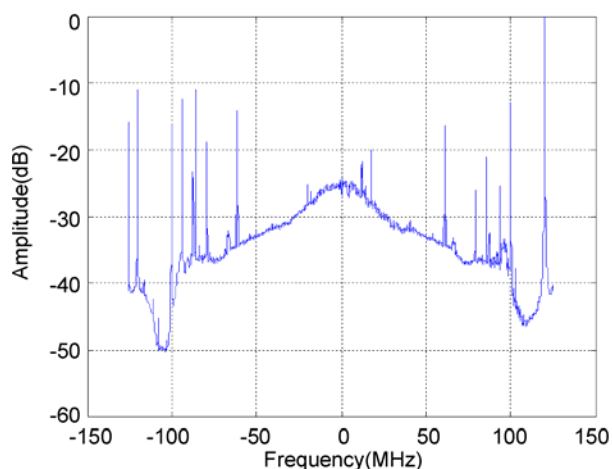


Figure 16. The measure spectrum of RFI.

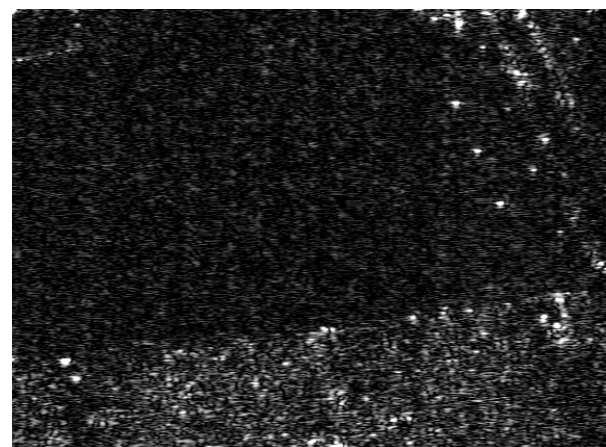
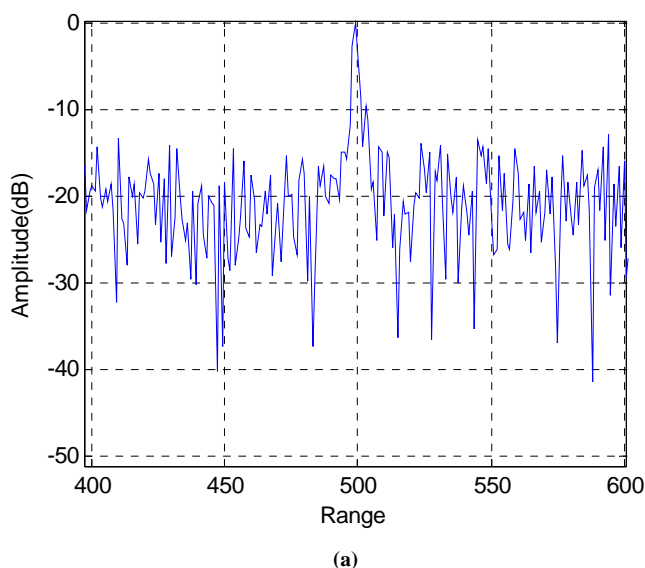


Figure 17. UWB SAR image without RFI suppression.

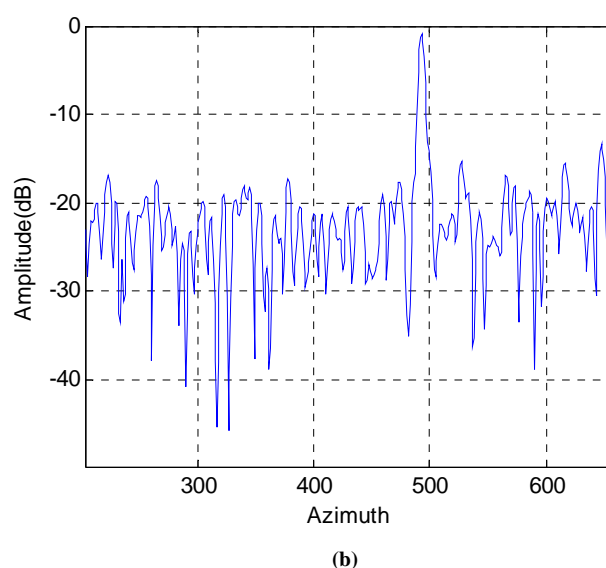


Figure 18. UWB SAR image without RFI suppression (a: The sectional plot in range, b: The sectional plot in azimuth).

NCS algorithm is applied to form the UWB SAR image. Figure 17 is a UWB SAR image without RFI suppression. Figure 18 is a cross-section of a point target corresponding to the frame in Figure 17, where Figure 18(a) is the plot along range and Figure 18(b) is the plot along azimuth.

Figure 19 is a UWB SAR image with RFI suppression adopting notch filter. Figure 20 is a cross-section of a point target in the frame of Figure 19, where Figure 20(a) is the plot along range and Figure 20(b) is the plot along azimuth.

We can observe that Figure 17 is blurred with RFI, especially in the relatively dark areas where have low SNR. Meanwhile, Figure 19 has more spots around strong point targets. These spots might be considered targets. But in the dark areas of Figure 19, the contrast of the image clearly improves.



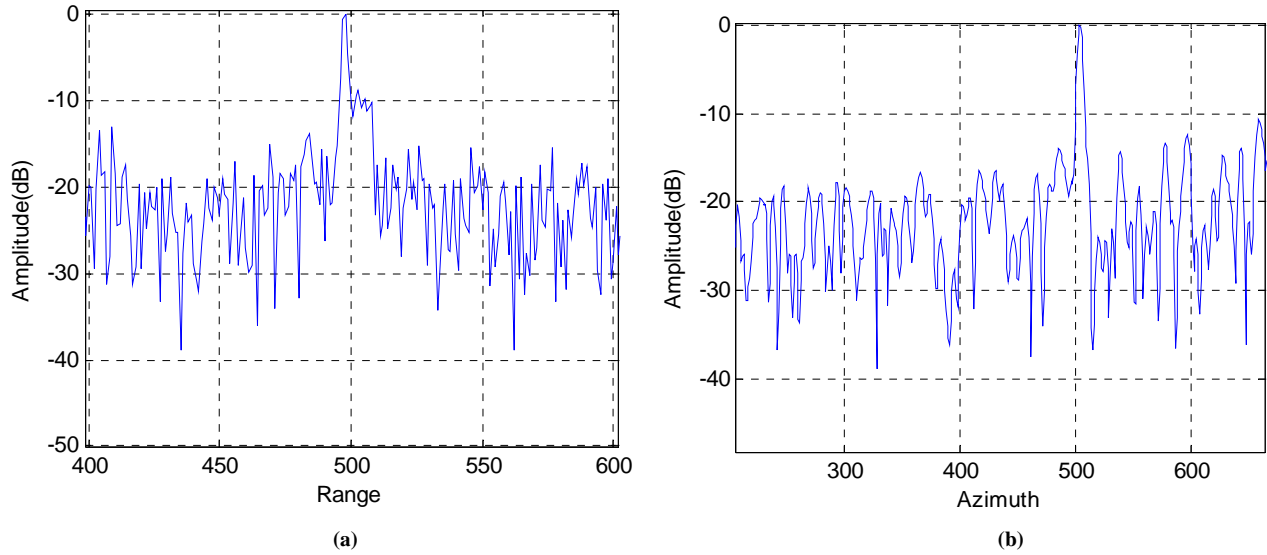


Figure 20. UWB SAR image without RFI suppression (a: The sectional plot in range, b: The sectional plot in azimuth).

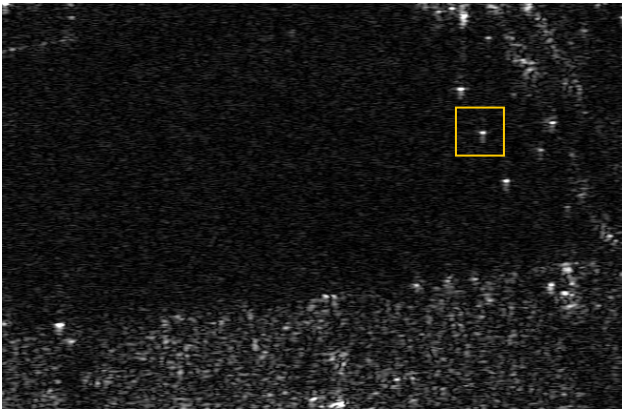


Figure 19. UWB SAR image with RFI suppression.

Table 3. Performance of the UWB SAR image.

	Resolution(m)	PSLR(dB)	ISLR(dB)
Without RFI suppression	1.2(r)×2.5(a)	−9.88	−4.79
With RFI suppression	1.2(r)×2.5(a)	−9.38	−2.98

\*Note: ISLR is calculated only in range.

The cross-section of a point target in Figure 18 and Figure 20 is used to analyse the influence before and after RFI suppression. Detailed results are described in Table 3. We find from Figure 18, Figure 19 and Table 3 that: 1) the average power of clutter, especially in range, is reduced after RFI suppression and 2) several spikes appear after RFI suppression adopting notch filter. Fortunately, the notch filter has little impact on the SAR resolution.

## 6. Conclusions

In this paper, the performance of the notch filter is theoretically analysed. Firstly, a matched filter with notch(s) is designed and its theoretical output is derived. The theoretical result shows that the performance of the notch filter is influenced by the notches' width and carrier, and that the notch filter has significant effect on the sidelobes but little effect on the resolution of system impulse response. Secondly, the simulation data and a specified system data are processed in range direction (one-dimension) with notch filter applied to test the validity of the theoretical result. Thirdly, the simulation data and the experimental UWB SAR data are applied to test the validity of the theoretical result. Some useful conclusions are made.

However, despite its shortages the notch filter is an effective tool to suppress RF interference for a small number of narrowband interferers in the lower SNR area of an image.

The theoretical result has been made available to describe the effect of the notch filter on the SAR image. We hope our work and its results may be helpful to who engaged in UWB SAR, wireless system design and RFI suppression research.

## 7. References

- [1] T. Koutsoudis and L. Lovas, "RF interference suppression in ultra wideband radar receivers," in *Algorithms for Synthetic Aperture Radar Imagery II* (D.A. Giglio, ed.), SPIE, Orlando, FL, Vol. 2487, pp. 107–118, April 1995.

- [2] T. Miller, L. Potter, and J. McCorkle, "Army research laboratory RFI suppression for ultra wideband Radar," *IEEE Transactions on Aerospace and Electronic Systems*, Vol. 33, No. 4, pp. 1142–1156, October 1997.
- [3] R. T. Lord and M. R. Inggs, "Efficient RFI suppression in SAR using LMS adaptive filter integrated with range/Doppler algorithm," *Electronics Letters*, Vol. 35 No. 8, pp. 629–630, April 15, 1999.
- [4] R. T. Lord and M. R. Inggs, "Approaches to RF interference suppression for VHF/UHF synthetic aperture radar," *Communications and Signal Processing, COMSIG'98*, pp. 95–100, 1998.
- [5] X. Luo, L. M. H. Ulander, J. Askne, G. Smith and P. O. Frolind, "RFI suppression in ultra-wideband SAR systems using LMS filters in frequency domain," *Electronics Letters*, Vol. 37, No. 4, February 15, 2007.
- [6] G. W. Davidson, I. G. Cumming, and M. R. ITO, "A chirp scaling approach for processing squint mode SAR data," *IEEE Transactions on Aerospace Electronic Systems*, Vol. 32, pp. 121–133, 1996.

# Reconfigure ZigBee Network Based on System Design

Yuan XU, Shubo QIU, Meng HOU

*School of Electronic Information and Control Engineering, Shandong Institute of Light Industry, Jinan, China*

*E-mail: qiushubo@163.com*

*Received March 13, 2009; revised May 12, 2009; accepted May 20, 2009*

## Abstract

This article analyses key technology used by network layer based on ZigBee technology. Then a reconfigure network as well as its strategy of forming network and distributing node is given. The simulation proved that the stability of reconfigure network and the ability of transmitting pass through obstacle are better than traditional network; it has an active significance for shorter delay because of the flexible of the improved forming network strategy.

**Keywords:** ZigBee, Reconfigure Network, Forming Network Strategy, Routing, NS-2

## 1. Introduction

ZigBee wireless communication technology which has the characteristics of short range, low power, low bit rate, and low cost is a new member of WLAN families. Actually, although it is used in more and more filed of produce, it causes some problems. When the environment of the detect area is changed, for instance, once the node died or the communication link is obstructed by some things for some reason, the link is broken. It will influences communication efficiency and the ability of transmitting through obstacle of ZigBee network.

In order to adapt to the environment changes of the detecting area, reconfigure ZigBee network is introduced. When the network node has problems (node death or communication path blocked), the reconfigure network can change its network structure, repair break link quickly, which improve the ability of transmitting through obstacle of ZigBee network [1].

## 2. ZigBee Network Layer

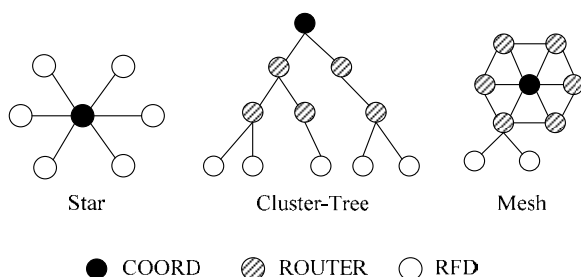
The main functions of ZigBee network layer are forming network, allocating address for node joined the network, routing discover, routing contain and so on. There are three types of topology as illustrated in Figure 1: Start, Cluster-Tree and Mesh network [2,3]. Based on the function of general network layer, ZigBee reduce energy consuming and link cost as much as possible.

In ZigBee networks, there are three types of nodes: 1) the ZigBee coordinator, which manages the network; 2) the routers, which are capable of participating in the AODV routing procedure; 3) the end devices, which transmit and receive frames through their parent node. Since end devices have no capability of AODV routing. Each non-ZigBee coordinator node has its parent while a non-end device node can have multiple children.

### 2.1. Forming Network

The first node which has coordination ability and has not joined any network coming into Personal Area Network (PAN) start to form network, and this node is the PAN Coordinator (PAN Coord) of this network. PAN Coord choose a idle channel after scanning channels, make sure its 16 bits network address, PAN ID, network topology parameters and so on. Then, it can accept other nodes as its children node.

When node A wants to join the PAN, it sends association request to the nodes in the network. If the node which receives the request has ability to accept node A to



**Figure 1. Network topology of ZigBee.**



be its children node, it assigns an unique 16 bits network address and sends association request to node A. Node A joins the network successfully and receives other nodes' association request after it receives association request. Whether node has the ability to receive and associate with other nodes is line to the source used by node, such as memory, energy and so on. If node in the network wants to leave network, it can sends remove association request to its parent node. It can leave network after receiving remove association request. What's more, the node must remove all the association with other nodes before leaving network if it has children nodes.

## 2.2. Network Address Allocate Mechanism

In ZigBee networks, the network addresses are assigned using a hierarchical addressing scheme. The address is unique within a particular network and is given by a parent to its children. The ZigBee coordinator determines  $C_m$  that is the maximum number of children a parent can have,  $R_m$  that is the maximum number of routers a parent can have as children, and  $L_m$  that is the maximum depth in the network.  $Cskip(d)$ , the size of the address sub-block being distributed by each parent at depth  $d$  to its router-capable child devices, is computed as follows [3].

$$Cskip(d) = \begin{cases} 1 + C_m \cdot (L_m - d - 1), & R_m = 1 \\ \frac{1 + C_m - R_m - C \cdot R_m^{L_m - d - 1}}{1 - R_m}, & \text{otherwise} \end{cases} \quad (1)$$

Network addresses  $A_{d+1, rn}$  and  $A_{d+1, el}$  will be assigned to the  $n$ -th router child and  $l$ -th end device child at depth  $d + 1$  in a sequential manner, respectively, according to the following equations:

$$A_{d+1, rn} = A_{parent} + Cskip(d) \times (n-1) + 1 \quad (2)$$

where  $1 \leq n \leq R_m$ , and  $A_{parent}$  represents the address of the parent.

$$A_{d+1, el} = A_{parent} + Cskip(d) \times R_m + 1 \quad (3)$$

## 2.3. ZigBee Routing Protocol

In order to achieve low cost and low power, ZigBee routing algorithm combines Cluster-Tree and Ad-hoc On-demand Distance Vector routing (AODV). But the AODV routing algorithm used by ZigBee is different from classical AODV routing algorithm, it should be called AODV Junior (AODVjr) routing algorithm accurately [4].

### 2.3.1. Cluster-Tree Routing Algorithm

In the Cluster-Tree Routing Algorithm, node calculates next hop according to network address of destination node. For example, the routing node which address is A

and depth is  $d$ , the node is its children node which address is  $D$  according to the following equations:

$$A < D < A + Cskip(d - 1) \quad (4)$$

If the destination node of group is the generation of the receiving node, node sends the group to its children node. At the time, set the address of next hop node is  $N$ , if  $D > A + R_m \times Cskip(d)$ ,  $N = D$ , otherwise,  $N$  is according to the following equations:

$$N = A + 1 + \left\lceil \frac{D - (A + 1)}{Cskip(d)} \right\rceil \times Cskip(d) \quad (5)$$

### 2.3.2. AODVjr Routing Algorithm

The AODV routing algorithm used by ZigBee is different from classical AODV routing algorithm, it should be called AODV Junior (AODVjr) routing algorithm accurately [4]. AODVjr has main function of AODV and make some simplifications for reducing cost, energy saving and so on.

1) AODVjr do not use sequence number in order to reduce control cost and simply routing progress. The sequence number used by AODV is to make sure that no cycle node at any time. AODVjr formulate that the destination node can reply RREP when it has group. Even some inter node having routing to destination node can not reply RREP.

2) AODVjr do not have precursor list, which simply routing table structure [5]. In AODV, if node detects there is an interruption in next-hop link, it makes upper node send RERR to inform influence node. AODVjr did not have precursor list because RERR is only sent to node which send failure.

3) AODVjr use local repair mechanism when link breaks. In the progress of repairing, it also dose not use sequence number but allow destination node to send RREP. Sending RERR to destination of data group and noticing that node no arrive if it repairs fail.

4) The node using AODV provides consistency information to other nodes by sending HELLO group periodically. The node using AODVjr dose not send HELLO group periodically, it updates neighbor list only according to receiving group or information provided by MAC.

### 2.3.3. ZigBee Routing

ZigBee Routing includes two types of node: RN+ and RN-. RN+ is a kind of node which has enough memory and be able to run AODVjr. RN- is contrary to RN+, it has neither enough memory nor ability of running AODVjr, so it has to process group by Cluster-Tree algorithm.

In Cluster-Tree algorithm, node delivers group to next hop immediately after receiving it. There is no routing progress and it is not necessary to contain routing table entry, so it reduces link cost and node energy consuming, what's more, it also reduce the demand to the memorial

ability of a node. However, because the link which built by Cluster-Tree may not be the best optimize link, the delay of deliver is big. The lower the depth of the node is, the more node works, and it makes assignment of communication flow unbalance. ZigBee routing allow RN+ to seek the best optimize link: RN+ start to use routing discover to find the shortest link to the destination node after it receives group, if there are two links have the same hop, node choose the better one according to the LQI which provides by MAC layer of IEEE802.15.4; node deliver group via the link, if there is an interruption in the link, RN+ repair routing by local repair mechanism. AODVjr reduce the delay of group deliver and improve reliability of group deliver [6].

### 3. Reconfigure ZigBee Network

#### 3.1. Reconfigure Network and Routing Algorithm

Reconfigure ZigBee Network is a network that the network structure can change dynamic. There are two meanings of structure changing dynamic: 1) the network topology does not change, but the type of node or the depth of network changes, 2) the network topology is changed.

According to above, in ZigBee technology, Cluster-Tree algorithm is suit for the node which the memory is limited, node can send the receiving group to its children node or parent node without routing discovery progress, but it have a big per-to-per delay in the network which distribute nodes unbalance, it can causes business flow distribute unbalance. The less the depth of the node is, the more node works, it makes low depth node exhaust its power more quickly and lead communication to break at last. In the contrary, ADOVjr algorithm can seek the best optimize link, reduce per-to-per delay, and relax business flow distribute unbalance. But ADOVjr need more memory to contain routing table and much link cost.

Reconfigure network use Cluster-Tree + AODVjr as its routing protocol, which can seek the best optimize link, reduce per-to-per delay, reduce link cost, and relax business flow distribute unbalance.

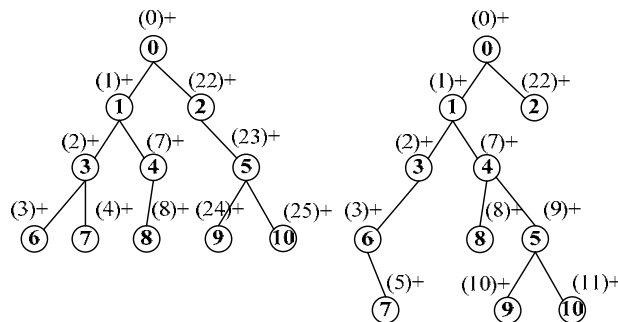
#### 3.2. Network Address Allocate Mechanism in Reconfigure Network

The reconfigure network use RN+ in order to reconfigure network to adapt the changing environment. The node joined the network form Cluster-Tree by the association provided by MAC layer. When a node allows a new node to join the network, they form parent-children relation; parent node assigns the unique 16 bits network address to children node. The Start topology is descript as a special Cluster-Tree topology (depth is 2) in reconfigure network.

#### 3.3. Network Reconfiguration

The network reconfiguration is embodied in the condition that the node environment is changing. Node need to reconfigure its network structure to communicate through obstruct when detect environment is changed.

As shown in Figure 2(a), as a result of changing environment, the communication link between node 7 and node 3 is obstructed in a Cluster-Tree network. Node has to leave network, it can't transmit data to upper lever node. At this time, node starts to reconfigure network: according to network address allocate mechanism in reconfigure network (Equations (1) (2)), compute network address of the node (node 6) with same depth and parent node. Then AODVjr routing algorithm is used for routing for this node. If find node 6, node 7 will join the network as children node of node 6 and re-allocate network address. The reconfiguration is finished when the communication link between node 6 and node 7 is built. If there is no node that both depth and parent node are the same around obstructed node (e.g. node 5), node will



**Figure 2. The forming network strategy of reconfigure Cluster-Tree network.** “—”is represent parent-children relation; “()” is represent network address; “+” is represent node RN+, “-” is represent node RN-; Node 0 is PAN Coord, its network address is 0.

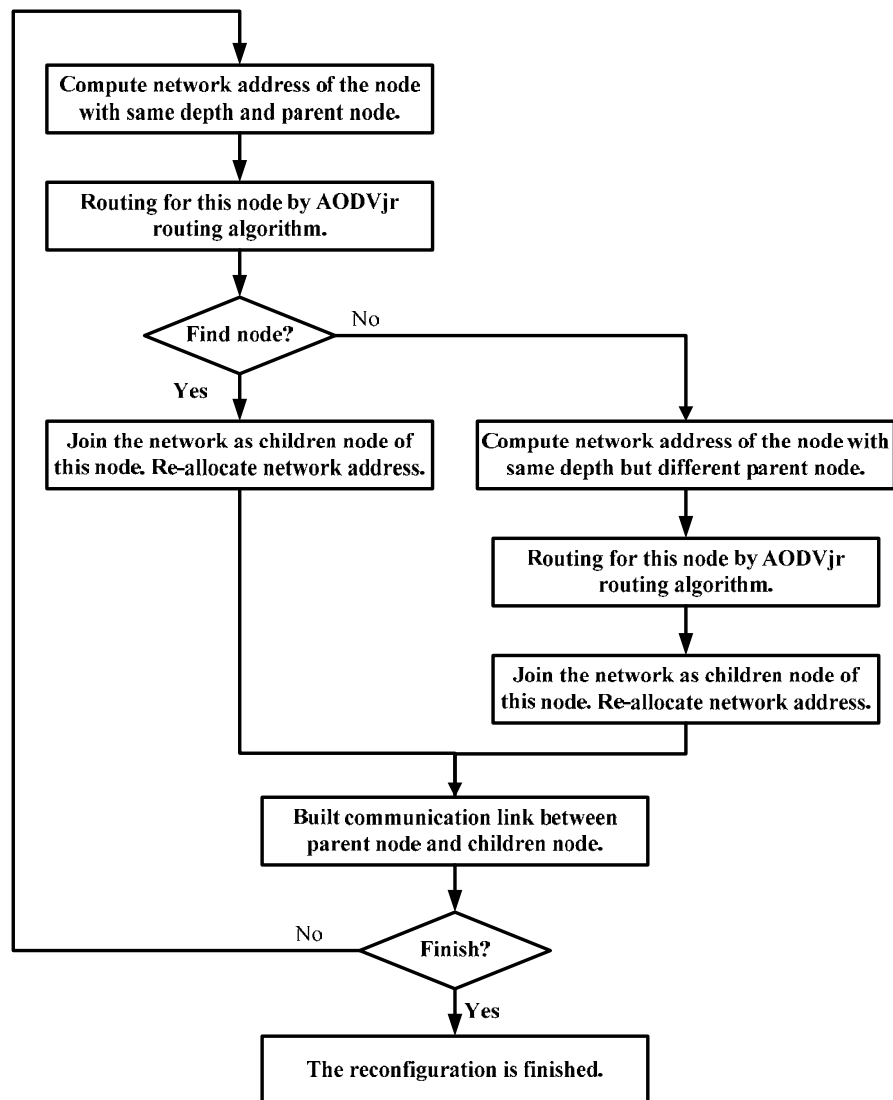


Figure 3. The process of reconfigure network.

compute network address of the node (node 4) with same depth but different parent node. Then according to network address allocate mechanism in reconfigure network, AODVjr routing algorithm is used for routing for this node. Node 5 will join the network as children node of node 4 and re-allocate network address if find node 4. The reconfiguration is finished when the communication link between node 4 and node 5 is built. Though the network structure is also Cluster-Tree after network configuration, the depth of network is changing from 4 to 5, as shown in Figure 2(b). Node 7's parent node is changing from node 3 to node 6 while node 5's parent node is changing from node 2 to node 4. The process of reconfigure network is shown in Figure 3.

The Start topology is defined as a special Cluster-Tree topology (depth is 2) in reconfigure network. So as

shown in Figure 4(a), according to AODVjr routing algorithm, node begins routing node 5 when the communication between node 4 and coordinator. Node 4 will join the network as children node of node 5 if find node 5 and be re-allocated network address. The reconfiguration is finished when the communication link between node 4 and node 5. As shown in Figure 4(c), the network topology is changing from Start to Cluster-Tree.

The value of network depth is limited. It is consider that the communication link is break and need to repair as long as network depth is over limited value.

### 3.4. Communication Link Repair

Routing across break node is a way of addressing which

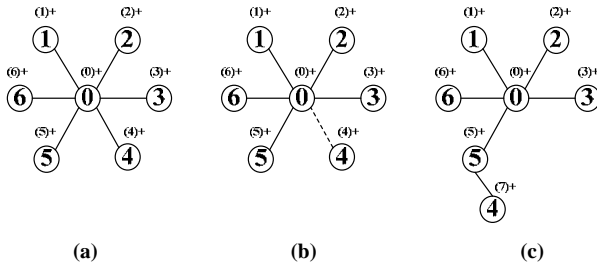


Figure 4. The forming network strategy of reconfigure start network.

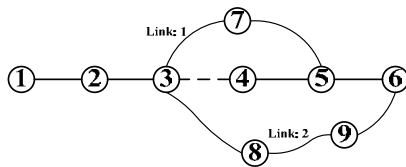


Figure 5. Routing across break node based on AODVjr. Link 1: If routing discovery find node 5, the new communication link will be 1-2-3-7-5-6. Link 2: If routing discovery can not find node 5, the new communication link will be 1-2-3-7-5-6.

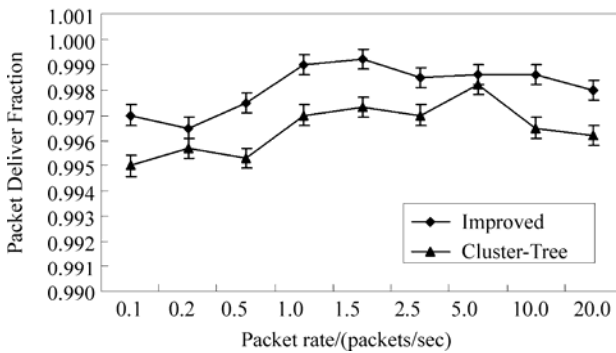
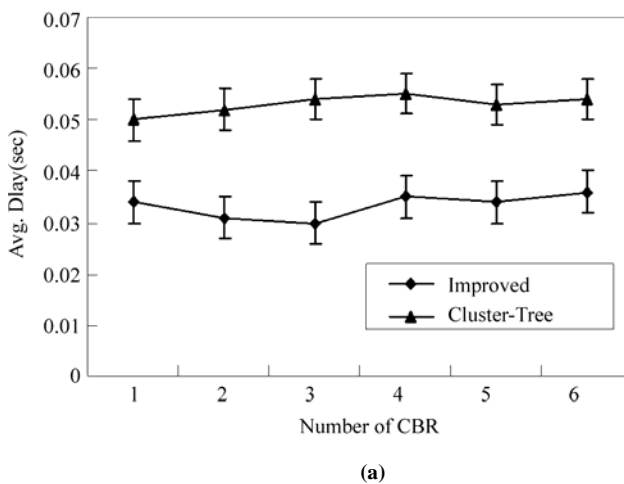
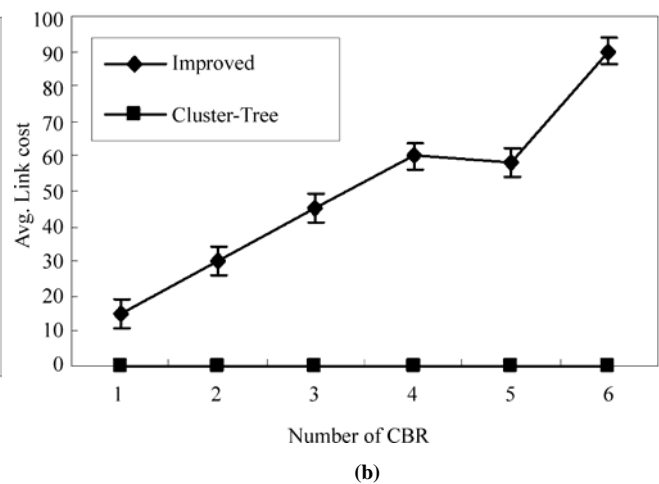


Figure 6. Performance with varying number of packet rate.



(a)



(b)

Figure 7. Performance with varying number of CBR.

seeks next two-hop node skip next node. In Figure 5, if node 4 fail in a five-hop distance communications link, node 3 can not receive MAC ACK. Considering that the communication links disconnected, node 3 will send a RERR packet to the node waiting for RREQ packet and set the path to invalid. Then, node 3 will find the address of next two jump node (node 5) and re-launch routing discovery to search for node 5. If routing discovery find node 5, the local repair will be completed by replacing the broken part of the communication link with a new rebuilding communications link (via node 7). The new communication link is 1-2-3-7-5-6 now. On the contrary, if routing can not find node 5, node 3 will find the address of destination node (node 6) and re-launch routing discovery to search for it. At last, the communication link between node 3 and node 6 is rebuilt. At the same time local repair is completed.

#### 4. Simulation Results

We used NS-2 Ver 2.33 as a simulator for performance evaluation of the reconfigure network, because the principal part of reconfigure network strategy is Cluster-Tree topology, we only simulate Cluster-Tree topology as shown in Figure 2. Parameters used in simulations are as following. Traffic sources are CBR (continuous bit-rate), MAC protocol is IEEE802.15.4, simulate time is 500 simulation seconds, data packet is 70 bytes; the entire node is RN+. Parameters of Cluster-Tree are:  $R_m=4$ ,  $C_m=4$ ,  $L_m=5$ . Both reconfigure network and Cluster-Tree network's packet deliver fraction in different deliver rate are shown in Figure 6, the average delay in different CBR is shown in Figure 7(a), and the link cost in different CBR is shown in Figure 7(b).

As shown in Figure 6, we can see packet deliver fraction of reconfigure network is higher than Cluster-Tree with the deliver rate growing, there are two reasons: 1) The Cluster-Tree topology transmit group by the parent-children relation. The lower the depth of the node is, the more data groups the node has to be progressed. The possibility of collision in MAC layer is so high that the failure of packet deliver is high. But in reconfigure network, lower depth node's works is reduced because the use of RN+ node. 2) There is only one link between source node and destination node, once the link breaks, it makes some groups can not be transmitted to destination node. Reconfigure network use Cluster-Tree + AODVjr algorithm, if link breaks, local repair can be used, which improved deliver successful.

As shown in Figure 7(a), the average per-to-per delay of reconfigure network is below 0.04s with the number of CBR increasing, while Cluster-Tree is higher than 0.05s. Link cost of reconfigure network is increasing along with the number of CBR increasing as shown in Figure 7(b).

## 5. Conclusions

ZigBee wireless communication technology has the characteristics of short range, low power, low bit rate, low cost; it has broad application prospects and enormous commercial value. This article analyses key technology used by network layer in detail. Then a reconfig-

ure network is given. The simulation proved that the stability of reconfigure network and the ability of transmitting through obstacle are better than traditional network; it has an active significance for shortening delay because of the flexible of the improved forming network strategy.

## 6. References

- [1] IEEE Standards 802.15.4TM-2003, Wireless medium access control (MAC) and physical layer (PHY) specifications for low-rate wireless personal area networks (LR-WPANs) [S].
- [2] ZigBee Alliance, ZigBee Specification, ZigBee Document 053474r06 Version 1.0 [S], 2004.
- [3] ZigBee Alliance, ZigBee Specification Version 2.0, December 2006, <http://www.zigbee.org>.
- [4] C. Perkins, E. Belding-Royer, and S. Das, RFC 3561, "Ad hoc on demand distance vector (AODV) Routing [S]," July 2003.
- [5] I. D. Chakeres and K. B. Luke, AODVjr, AODV simplified [J], ACM SIG-MOBILE Mobile Computing and Communications Review, Vol. 6, No. 3, pp. 100–101, 2002.
- [6] L. Kleinrock and F. A. Tobagi, "Packet switching in radio channels: Part 1-Carrier sense multiple-access modes and their throughput-delay characteristics," IEEE Transactions on Communications, Vol. COM-23, pp. 1400–1416, December 1975.

# Optimal Deployment with Self-Healing Movement Algorithm for Particular Region in Wireless Sensor Network

Fan ZHU, Hongli LIU, Shugang LIU, Jie ZHAN

*College of Electrical and Information Engineering, Hunan University, Changsha, China*

*E-mail: henry1214ars@yahoo.com.cn*

*Received July 29, 2009; revised August 15, 2009; accepted August 20, 2009*

## Abstract

Optimizing deployment of sensors with self-healing ability is an efficient way to solve the problems of coverage, connectivity and the dead nodes in WSNs. This work discusses the particular relationship between the monitoring range and the communication range, and proposes an optimal deployment with self-healing movement algorithm for closed or semi-closed area with irregular shape, which can not only satisfy both coverage and connectivity by using as few nodes as possible, but also compensate the failure of nodes by mobility in WSNs. We compute the maximum efficient range of several neighbor sensors based on the different relationships between monitoring range and communication range with consideration of the complex boundary or obstacles in the region, and combine it with the Euclidean Minimum Spanning Tree (EMST) algorithm to ensure the coverage and communication of Region of Interest (ROI). Besides, we calculate the location of dead nodes by Geometry Algorithm, and move the higher priority nodes to replace them by another Improved Virtual Force Algorithm (IVFA). Eventually, simulation results based-on MATLAB are presented, which do show that this optimal deployment with self-healing movement algorithm can ensure the coverage and communication of an entire region by requiring the least number of nodes and effectively compensate the loss of the networks.

**Keywords:** Optimal Deployment, Self-Healing Movement, Particular Region, Euclidean Minimum Spanning Tree (EMST), Improved Virtual Force Algorithm (IVFA)

## 1. Introduction

Sensor networks consist of a large number of small, light-weight, highly power-constrained, and inexpensive wireless nodes called sensors. Sensors are equipped with detectors for intrusion, sensing changes in temperature, humidity, chemicals, or any other characteristic of the environment that needs to be monitored. The data about the environment is constantly observed, consolidated, and sent to a monitor or Base Station (BS). Data transmission from the sensors to the BS can be periodic, event-triggered, or in response to a query from the BS. While each sensor node has limited computation capabilities and usually non-rechargeable battery power, the collaboration among thousands of sensors deployed in a region makes sensor networks a powerful system for observation of the environment [1]. The data sensed by

the sensors is generally highly critical, and may be of scientific or strategic importance. Hence, the coverage provided by sensor networks is a very important criterion of their effectiveness. Special emphasis is placed on coverage especially in tactical applications such as surveillance and reconnaissance. Sensors can easily be used for the perilous and demanding duties of observing landscapes for intrusion detection [2]. In a wireless sensor network, the reasonable deployment of sensors should take both coverage and connectivity into account. Coverage requires that any physical field in a sensing region can be monitored by at least one node. Connectivity requires that each node is under the range of communication of its neighbor sensors. All these nodes can consist of an Ad-hoc network, and also transmit data packets to the BS. On the other hand, as time progresses the sensor nodes may die randomly due to malfunction, energy ex-

haustion or malicious destruction. All these factors result in holes of coverage and connectivity in WSNs, which makes the system unable to meet the performance criterion.

## 2. Related Works

Recently, a lot of approaches were advanced to solve these problems of coverage, connectivity and dead nodes in WSN. The work in [3,4] discusses how to adjust the locations of the nodes to satisfy the coverage in an open space, but without considering the boundary or obstacle. The grid algorithm in [5,6] is an appropriate way to ensure coverage and connectivity when there are a few nodes in the sensing region, however, with increasing nodes, it would be low efficient. The work in [7] does consider both coverage and communication, but it defaults that the range of sensing and communication are equal to each nodes without discussing varied situations respectively because of the different ranges between sensing and communication. The intersection point of the two dead nodes' neighboring sensors is used to decide where available node moves towards based on the random deployment in [2,8]. The Virtual Force Algorithm (VFA) strategy to enhance the coverage after an initial random placement of sensors is proposed by [9,10], but the shortest distance path among nodes is not rectified during the process of movement. Works [11,12] efficiently adjust the sensor placement after an initial random deployment and apply fuzzy logic theory to handle the uncertainty in sensor deployment problem.

## 3. Problem Formulation

Assuming a sensing field, the range of communication of each sensor in this region is  $R_c$ , within which it can transmit data packets to other sensors. Also, the sensing distance is  $R_s$ . The areas of each node's coverage and connectivity are assumed as two ideal circles respectively. K. Kar and S. Banerjee default  $R_c=R_s$  in work [7], which satisfy the coverage and connectivity of the region. However, it is not realistic to analyze this issue by defining  $R_c=R_s$  simply. We discuss different situations based on different relationships between  $R_c$  and  $R_s$  in order to adopt an appropriate deployment approach. According to literature [13] D. Pompili the two adjacent sensors, which are separated by no more than  $\sqrt{3} R_s$ , can ensure effective coverage of the surrounding region. Thus, we can figure out the deployment approaches that respectively regard coverage or connectivity as the first choice in open space. Our reform does not let nodes restrict by any choice, but simultaneity satisfy coverage and connectivity by the least number of nodes. Secondly, another important issue is how to deploy sensors effectively

in that region with boundary and obstacle. Because the boundary or the obstacle may limit the distance of sensing and communication, the approach of placing this kind of region is not different from placing in open space. On the other hand, the sensor nodes may die randomly due to malfunction, energy exhaustion or malicious destruction as time progresses. All these factors result in holes of coverage and connectivity in WSN, which makes the system unable to meet the performance criterion. In this paper, we also propose another movement approach to compensate the loss based on the mentioned optimal deployment algorithm. Below, we discuss how to deploy in particular region.

## 4. Optimal Deployment with Self-Healing Movement Algorithm for 'Particular Region'

The 'Particular Region' is a closed area or semi-enclosed area with boundary or obstacle which is consisted of unregulated polygons and arches. The optimal deployment algorithm can ensure this region's coverage and connectivity by the least sensors. However, researching on the deployment algorithm in an open space is regarded as the foundation to analyze the different situations in a particular region. This work discusses the deployment in an open space firstly, and then we improve and summarize the specific deployment algorithm in a particular region. All nodes are deployed above the ground about one meter to ensure the most optimal channel.

### 4.1. Deploying Sensors in an Open Space

Multi-line sensor arrays effectively resolve the issue. Below we study on deploying sensors in an open area without obstacles, and then extend to the deployment method in particular area with boundaries and obstacles.

Firstly, we establish a two-dimensional coordinates without boundaries or obstacles, and deploy lines of sensors, it guarantees the entire coverage of both adjacent nodes and each row. As the adjacent nodes can communicate with each other, if it is required to maintain the whole region's connectivity, we can add some of the sensors between adjacent lines to ensure it.

*Case 1:*  $R_c \leq \sqrt{3} R_s$ , the distance of adjacent nodes at each line is set  $R_c$ , which guarantees the coverage of adjacent nodes. Because  $R_c \leq \sqrt{3} R_s$ , the width of belt-like region that covered by a row of sensors is

$$2 \times \sqrt{R_s^2 - \frac{R_c^2}{4}}.$$

The difference value of two adjacent lines of nodes on the Y-axis is  $R_s + \sqrt{R_s^2 - \frac{R_c^2}{4}}$ , while  $\pm \frac{R_c}{2}$  on the

X-axis. The above-mentioned method guarantees the coverage of the entire region. Figures 1–3 show three possible conditions. For  $R_c < \sqrt{3} R_s$ , what should be paid more attention to is that the method only satisfies communication property between adjacent nodes but not adjacent lines.

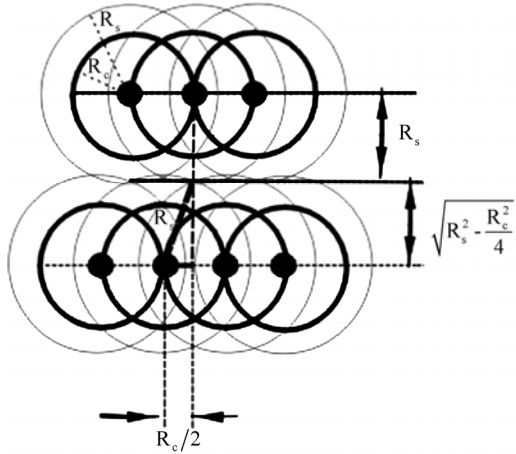


Figure 1. Deployment when  $R_s > R_c$ .

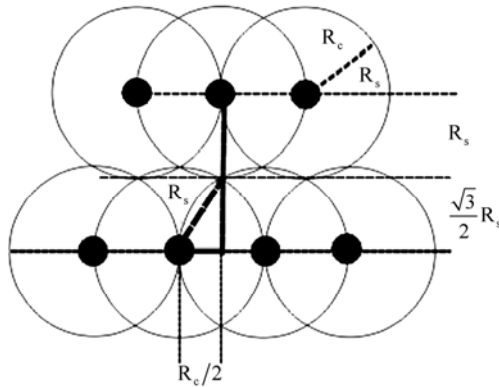


Figure 2. Deployment when  $R_s = R_c$ .

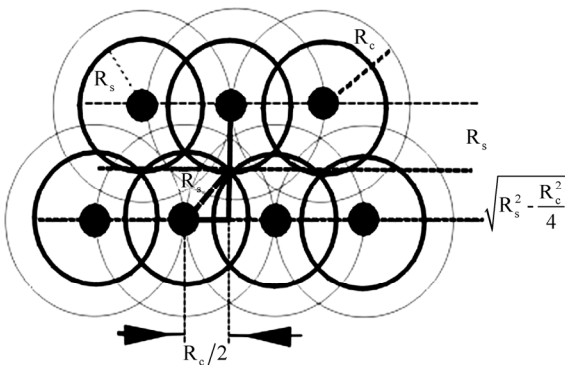


Figure 3. Deployment when  $R_s < R_c < \sqrt{3} R_s$ .

*Case 2:*  $R_c > \sqrt{3} R_s$ , as the smaller  $R_s$ , if we continue to adopt above-mentioned methods, which would lead the blind region that is not monitored between two adjacent lines, and also result in a waste of sensor nodes. So here's a typical use of the principle of hexagonal fabric which is more reasonable, and set the two adjacent sensors by  $R_s$  as Figure 4 shows. Therefore, it ensures the regional coverage and connectivity.

## 4.2. Deploying Sensors in Particular Area

For placing in particular regions, we can sum up the rules of deployment in two-dimensional coordinates from analyzing on a large area. First, establishing the two-dimensional coordinates, and assuming a initial node  $S_{(0,0)} = (x_0, y_0)$  which is nearest to the origin than other nodes. According to above conclusion, the node  $S_{(1,0)} = (x_0 + R_c, y_0)$ , which is deployed next to the initial node along the positive x-axis direction. While the node  $S_{(0,1)} = (x_0 + R_c/2, y_0 + R_s + \sqrt{R_s^2 - \frac{R_c^2}{4}})$  deployed first on the second row that is close to initial row along the positive y-axis direction.

$$S_{(2,2)} = (x_0 + 2R_c, y_0 + 2R_s + 2\sqrt{R_s^2 - \frac{R_c^2}{4}}).$$

By the same token any node's position placed by the above deployment algorithm in the two-dimensional coordinates can be calculated. (Note: The odd and even lines are different):

$$S_{(n', 2n)} = (x_0 + n'R_c, y_0 + 2nR_s + 2n\sqrt{R_s^2 - \frac{R_c^2}{4}}), \quad (1)$$

$$S_{(n', 2n+1)} = (x_0 + R_c/2 + n'R_c, y_0 + (2n+1)R_s + (2n+1)\sqrt{R_s^2 - \frac{R_c^2}{4}}), \quad (2)$$

$$n' = (0, 1, 2, \dots, \infty), n = (0, 1, 2, \dots, \infty).$$

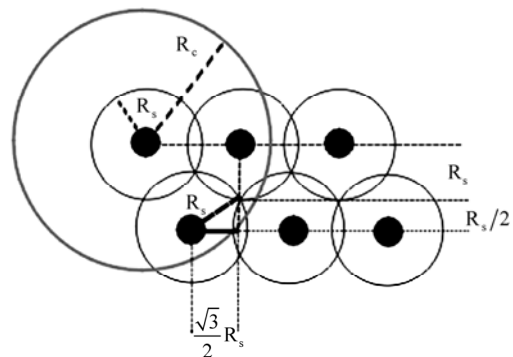


Figure 4. Deployment when  $R_c > \sqrt{3} R_s$ .



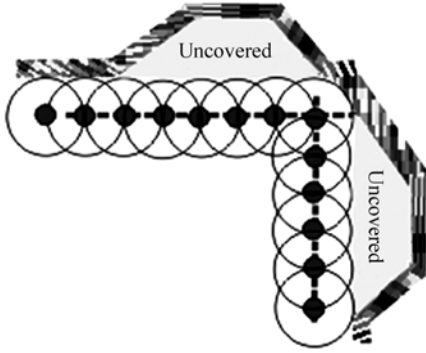


Figure 5. This approach leads to uncovered area.

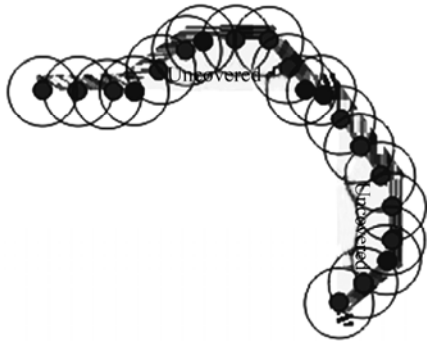


Figure 6. Deploying along the boundary (nodes in the corners are redundant).

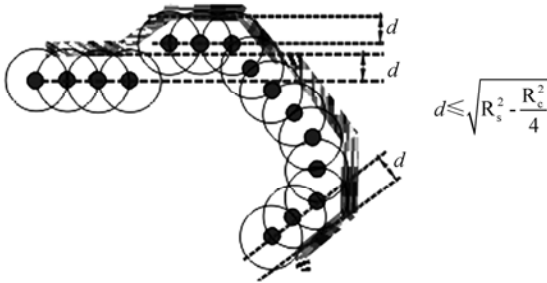


Figure 7. Improved deployment.

Of course, the node expression derived would change if we set different initial nodes. We need to be flexible in establishing the appropriate two-dimensional coordinates depending on the different shape of the area. However, the above-mentioned connectivity, which is limited that only ensures sensing data exchange between adjacent nodes at the same line, is limited. The following work focusing on the deployment of particular area will discuss how to ensure entire network's connectivity.

Assuming several nodes in a particular area with boundaries and obstacles are deployed as Figure 5, it results in uncovered area. The approach as Figure 6 shows meets the whole coverage and connectivity, but in

order to ensure the communication of adjacent nodes not to be blocked by obstacles, there are extra nodes to be added at the corner of the boundary, so it certainly wastes sensors. The deployment method in particular area can be improved on the basis of the above research as Figure 7 shows.

Assuming  $d$  is the width of uncovered area:

Case 1:  $d \leq \sqrt{R_s^2 - \frac{R_c^2}{4}}$ , and  $R_c \leq \sqrt{3} R_s$ , the distance

between nodes and boundary is set  $\sqrt{R_s^2 - \frac{R_c^2}{4}}$ , and two adjacent nodes are separated by  $R_c$ .

Case 2:  $d \leq \sqrt{R_s^2 - \frac{R_c^2}{4}}$ , and  $R_c > \sqrt{3} R_s$ , the distance

between nodes and boundary is set  $\sqrt{R_s^2 - \frac{R_c^2}{4}}$ , and two adjacent nodes are separated by  $\sqrt{3} R_s$ . This deployment method can satisfy coverage and connectivity.

Case 3:  $d > \sqrt{R_s^2 - \frac{R_c^2}{4}}$ , no matter what relationship

between  $R_c$  and  $R_s$  is, the method is as well as the deployment approach in large area.

In this way, both the coverage of whole area and the connectivity of adjacent nodes are guaranteed by the least number of nodes.

However, only satisfying connectivity of adjacent nodes on the same line is not enough to make all the nodes form an Ad-hoc network. In this paper, the EMST (Euclidean Minimum Spanning Tree) algorithm is introduced to estimate the communication location of the longest boundary, and also it combines with geometric analysis to solve the entire network problem of connectivity.

Assuming  $T$  area, and Choose any point  $S$  that can correspond to a leaf node of  $T$ , and all  $\{S\}$  are defined as subsets to  $C$  in  $T$ , set  $C \leftarrow \{S\}$ , Also set  $K=0$ ,  $K \rightarrow K+1$ . Choose any  $S' \in C$ . The  $R_c$ -disk which is chosen as  $D_k$  centered at  $S'$ . Move any points in  $C$  which are covered by  $D_k$ . Set  $I_k$  as the point of intersection by  $D_k$  and the boundary of  $T$ . For each point  $S'' \in I_k$ , including  $S'' \in C$ , if  $S'' \notin D_1 \cup D_2 \cup D_3 \cup \dots \cup D_{k-1}$ . So the path from initial  $S$  to  $S''$  in  $T$  is covered by  $D_1 \cup D_2 \cup D_3 \cup \dots \cup D_{k-1} \cup D_k$  completely. The deployment of specific path can be regarded as the geometric issues. The straight-line distance between two lines of

nodes is  $R_s + \sqrt{R_s^2 - \frac{R_c^2}{4}}$ . According to the parallelogram principle, two diagonal  $d_1, d_2$  respectively is:

$$d_1 = (R_s + \sqrt{R_s^2 - \frac{R_c^2}{4}})^2 + \frac{R_c^2}{4}; \quad (3)$$

$$d_2 = (R_s + \sqrt{R_s^2 - \frac{R_c^2}{4}})^2 + \frac{9R_c^2}{4}; \quad (4)$$

We choose an appropriate diagonal depending on the different shape of particular area, and calculate the number of complementary sensors:  $\frac{d_1}{R_c}$  or  $\frac{d_2}{R_c}$ . These added sensors are separated by the equal distance on the diagonal as shown in Figure 8. To sum up, the communication path among all lines is established, and all sensors ensure the connectivity in entire network.

### 4.3. Self-Healing Movement

Even if the applications of above-mentioned optimal deployment algorithm can satisfy both connectivity and coverage, the sensor nodes may die randomly due to malfunction, energy exhaustion or malicious destruction as time progresses. All these factors result in holes of coverage and connectivity in WSN, which makes the system unable to meet the performance criterion. In this paper, we propose another movement approach to compensate the loss based on the mentioned optimal deployment algorithm. After the initialization of network, we assume that every node is equipped with the capability of movement, and acquires their location and communication neighbors respectively by localization protocol as [14,15] referred. Also communication neighbors will detect when any node dies. Then these nodes' neighbors broadcast a packet containing its location to next one-hop node which continues to transmit to another until all the nodes get the message of hole. The following section presents our movement algorithm:

According to above-mentioned deployment algorithm, in order to satisfy the whole connectivity and coverage in networks, it is inevitable to produce some edge nodes

whose real coverage areas are smaller than other's like the R nodes as the Figure 9 shows. In our approach, we need to make full use of these edge nodes to compensate the holes of coverage and connectivity in WSN. Hence we divide those nodes which were deployed near to the boundary or obstacles into three categories on the basis of different relationship between  $R_s$  and  $R_c$ :

*Case 1: When  $R_c \leq \sqrt{3} R_s$ :*

1) The vertical distance between node and boundary or obstacles is  $d < \sqrt{R_s^2 - \frac{R_c^2}{4}}$ ,

2) The vertical distance between node and boundary or obstacles is  $\sqrt{R_s^2 - \frac{R_c^2}{4}} \leq d < R_s$ ,

3) The vertical distance between node and boundary or obstacles is  $d \geq R_s$ .

Also three types of nodes are set in different priority classes to move. The nodes ( $d < \sqrt{R_s^2 - \frac{R_c^2}{4}}$ ) get the top priority, while the nodes ( $\sqrt{R_s^2 - \frac{R_c^2}{4}} \leq d < R_s$ ) is mid.

*Case 2: When  $R_c > \sqrt{3} R_s$ :*

1) The vertical distance between node and boundary or obstacles is  $d < \frac{R_s}{2}$ ,

2) The vertical distance between node and boundary or obstacles is  $\frac{R_s}{2} \leq d < R_s$ ,

3) The vertical distance between node and boundary or obstacles is  $d \geq R_s$ .

Also three types of nodes are set in different priority classes to move. The nodes ( $d < \frac{R_s}{2}$ ) get the top priority, while the nodes ( $\frac{R_s}{2} \leq d < R_s$ ) is mid.

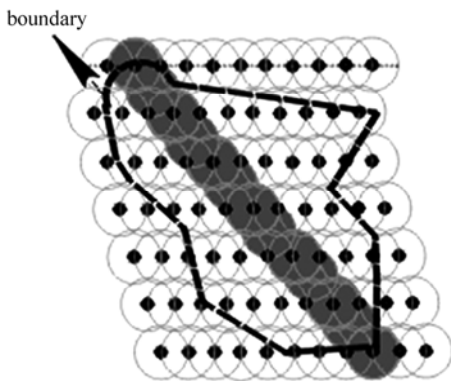


Figure 8. Deployment by EMST.

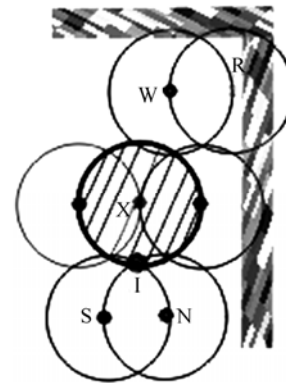


Figure 9. X is a dead node, R is the node with top priority.

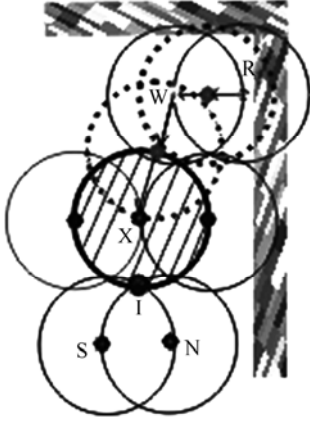


Figure 10. The process of movement.

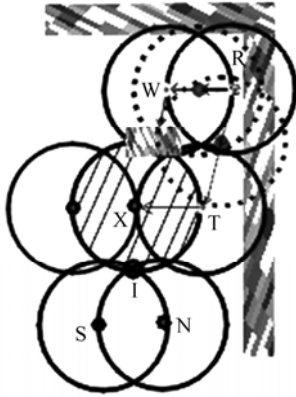


Figure 11. Changing another routing movement.

Though the following research and proof are based on  $R_c \leq \sqrt{3} R_s$ , we adopt the same principle methods to move the nodes on the condition of  $R_c > \sqrt{3} R_s$  except only the constant of slope and location of nodes changed by above analyses. When the message is received by the nodes with top priority, they can calculate the distance from the dead node to its location  $d_{dt} = \sqrt{(x_t - x_d)^2 + (y_t - y_d)^2}$ , (i.e.  $(x_t, y_t)$  ( $x_d, y_d$ ) are the location of top priority node and dead node in two-dimensional coordinates respectively) and broadcast to the other top priority nodes. The radio energy dissipation model as [16] referred:

$$E_T(l, d) = \begin{cases} lE_{elec} + l\varepsilon_{fs}d^2 & \text{if } d \leq 0 \\ lE_{elec} + l\varepsilon_{mp}d^4 & \text{if } d > 0 \end{cases} \quad (5)$$

It presents that more and more energy would be wasted with the increasing of distance. So we choose the minimum distance  $d_{dt}$  from all the top movement priority nodes. In order to avoid energy depletion caused by excessive movement, we propose the new type move-

ment like routing to counterpoise the node's energy consumption in movement by calculation of routing movement distance as shown in Figure 10. When the available node R with top priority is chosen to replace the dead node X, R will broadcast available nodes W, S, N and then these nodes produce the virtual force to make W move to the location of X, and simultaneously R is forced to the W location. On the other hand, if there is an obstacle caused failure of movement on the path from W to X, W will broadcast back to R, and the routing of movement change immediately as the Figure 11 shows.

#### 4.4. Calculation of Routing Movement Distance

Assuming the neighbor W of the dead node X. Set the co-ordinates of X be  $(x_0, y_0)$ , and those of W be  $(x_w, y_w)$ . Consider another two neighbors of node X, S and N which located at  $(x_s, y_s)$  and  $(x_n, y_n)$  respectively. The circle of coverage of nodes S and N intersect at the point I by the co-ordinates. Our algorithm makes node W move towards X such that the area that was earlier sensed by X can now be covered by node W.

*Step 1:* The co-ordinates of intersection node I and the distance  $d_{sn}$  between S and N can be derived as follows:

$$\begin{cases} (x - x_s)^2 + (y - y_s)^2 = R_s^2 \\ (x - x_n)^2 + (y - y_n)^2 = R_s^2 \end{cases} \quad (6)$$

$$d_{sn} = \sqrt{(x_s - x_n)^2 + (y_s - y_n)^2} \quad (7)$$

The solution of equation group:

$$\begin{cases} x = (x_s + x_n)/2 \pm (y_s - y_n)\sqrt{d^2(4r_c^2 - d^2)}/2d^2 \\ y = (y_s + y_n)/2 \mp (x_s - x_n)\sqrt{d^2(4r_c^2 - d^2)}/2d^2 \end{cases} \quad (8)$$

The solutions  $(x_i, y_i)$  which is closer to dead node X is the required answer.

*Step 2:* As above-mentioned optimal deployment algorithm, adjacent nodes have a fixed line slope in different relationship between  $R_s$  and  $R_c$ .

$$\text{When } R_c \leq \sqrt{3} R_s: \quad \text{tg} \alpha = (R_s + \sqrt{R_s^2 - \frac{R_c^2}{4}}) / \frac{R_c}{2}$$

$$\text{When } R_c > \sqrt{3} R_s: \quad \text{tg} \alpha = \sqrt{3}$$

*Step 3:* set  $X' = (x', y')$  as the point node W move towards.

$$\begin{cases} (x' - x_i)^2 + (y' - y_i)^2 = R_s^2 \\ (y_w - y') / (x_w - x') = \text{tg} \alpha \end{cases} \quad (9)$$

So we can prove that the node W move towards  $(x', y')$  which was the location of dead node X.

#### 4.5. Improved Virtual Forces Algorithm

In the process of movement, we also combined with the VFA model as [10] presents:

$$\vec{F}_{ij} = \begin{cases} (\omega_A(d_{ij} - d_{th}), \alpha_{ij}) & \text{if } d_{ij} > d_{th} \\ 0, & \text{if } d_{ij} = d_{th} \\ (\omega_R \frac{1}{d_{ij}}, \alpha_{ij} + \pi), & \text{if otherwise} \end{cases} \quad (10)$$

where  $d_{ij}$  is the Euclidean distance between sensor  $s_i$  and  $s_j$ ,  $d_{th}$  is the threshold on the distance between  $s_i$  and  $s_j$ ,  $\alpha_{ij}$  is the orientation (angle) of a line segment from  $s_i$  to  $s_j$ , and  $w_A(w_R)$  is a measure of the attractive (repulsive) force. The threshold distance  $d_{th}$  controls how close sensors get to each other. We assume that the neighbors of the dead nodes will produce the “attractive” force to the predetermined movement nodes. In order to reduce the total moving distance of the sensors, we determine whether  $s_i$  can move toward  $p_j$  at every period (namely round) as follows:

*Step 1:* The dead nodes  $p_j$  is detected by its neighbors, and its location is obtained by above geometry algorithm.

*Step 2:* Calculating  $d_{p_j s_1}, d_{p_j s_2}, \dots, d_{p_j s_n}$ . When the shortest  $d_{p_j s_i}$  is found, the  $s_i$  decide to move toward  $p_j$  with a threshold  $\lambda$ ,  $\lambda$  is the maximal distance a sensor can move forward at every round. Then the  $s_i(x_{s_i}, y_{s_i})$  is updated with  $s'_i(x'_{s_i}, y'_{s_i})$  which can be calculated by the Equations (1) and (2).

As Figure 12 shows, the linear equation of the line passes through the sensor  $s_i$  and the predetermined location  $p_j$  is  $(y - y_{p_j})(y_{s_i} - y_{p_j}) = (x - x_{p_j})(x_{s_i} - x_{p_j})$ . We can obtain

$$x'_{s_i} = \lambda(x_{s_i} + x_{p_j}) / d_{p_j s_i} + x_{p_j}$$

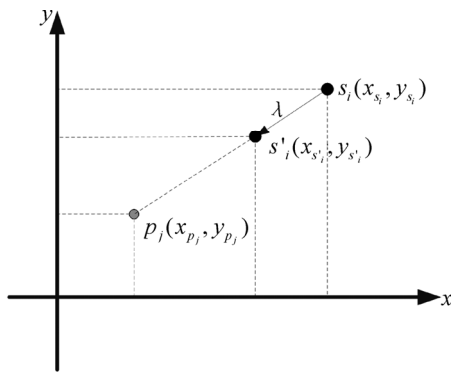


Figure 12. The coordinate of  $s_i$  is updated after moving a  $\lambda$  threshold.

and  $y'_{s_i} = \lambda(y_{s_i} + y_{p_j}) / d_{p_j s_i} + y_{p_j}$ .

So we can summarize an improved VFA with above analyses, if the final force of the dead node's neighbors is calculated, the sensor with priority moves towards the dead node's location according to the magnitude and direction. The updated move can be calculated:

$$\begin{cases} x(i)_{new} = x(i)_{old} + \text{sign}(\vec{F}_{ix}) \left| \frac{\vec{F}_{ix}}{\vec{F}_i} \right| \times \lambda \\ y(i)_{new} = y(i)_{old} + \text{sign}(\vec{F}_{iy}) \left| \frac{\vec{F}_{iy}}{\vec{F}_i} \right| \times \lambda \end{cases} \quad (11)$$

To sum up, with combination geometry and improved VFA, we prove the feasibility of our movement algorithm. If all these top priority nodes already compensate the loss in the network, the mid priority nodes will continue to move to the lower and dead ones. Note that the crucial Euclidean leaf nodes which ensure the whole connectivity of network need to be recovered first if any one doesn't work.

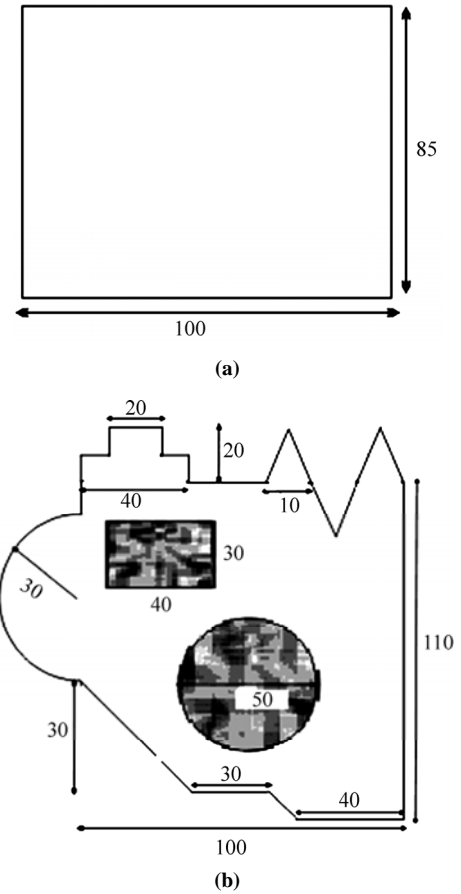
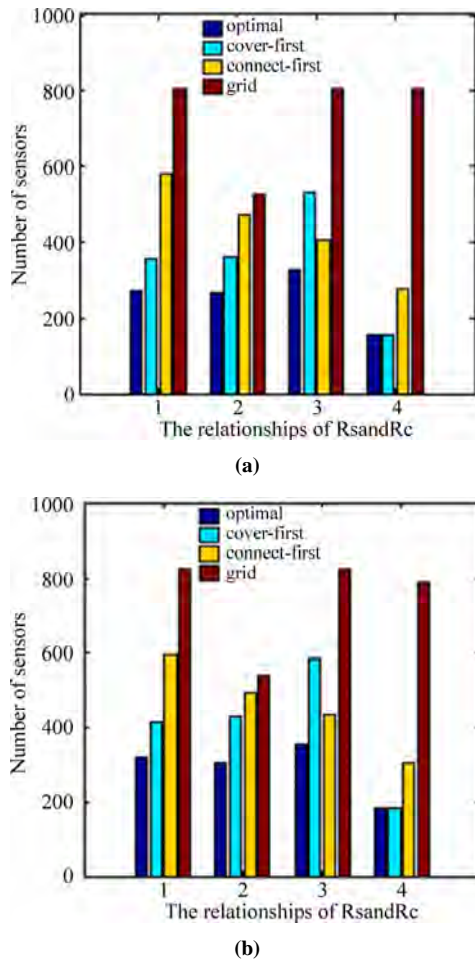


Figure 13. (a) 100\*85 rectangle, (b) Particular areas: a complex area with kinds of boundaries and obstacles (shadow).

## 5. Simulation Results

In this section, we present two groups of experimental results to prove the effectiveness of our optimal deployment with self-healing movement algorithm in different fields. We choose one simple and representative sensing fields, and then design a more complex particular area as shown in Figures 13(a), (b). We consider four groups of cases  $(R_c, R_s) = (4, 6); (5, 5); (6, 4); (8, 4)$  to reflect the relationships as above-mentioned:  $R_s > R_c$ ;  $R_s = R_c$ ;  $R_s < R_c < \sqrt{3} R_s$ ;  $R_c > \sqrt{3} R_s$  respectively. All nodes are deployed above the ground about one meter to ensure the most optimal channel. We compare the number of sensor being deployed as comparison metric in four different methods including ours optimal algorithm, coverage-first algorithm, connect-first algorithm, grid algorithm discussed in Section 3. Then we make some nodes die randomly on purpose, and compare the coverage of this network with the other one which already healed it.



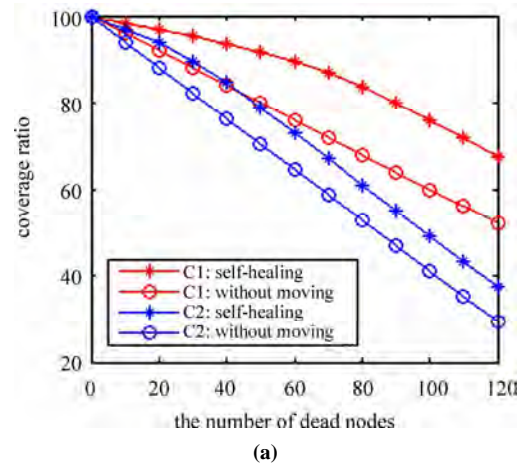
**Figure 14.** (a) The number of sensors used in 100\*85 rectangle, (b) The number of sensors used in particular sensing area.

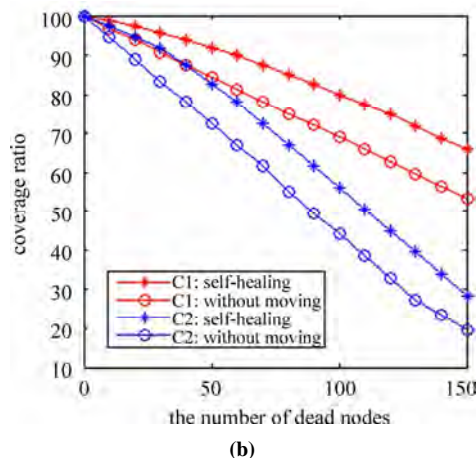
The number of sensors in four different relationships of  $R_s$  and  $R_c$  under particular area is shown in Figure 14. Most sensors are used in the grid algorithm, because all adjacent nodes are separated by the minimum of  $R_s$  and  $R_c$ . Sensors are placed horizontally by separation of  $R_c$  under the connect-first algorithm, thus it results in wasting many overlapping nodes in coverage when  $R_s > R_c$ . On the contrary, for it needs so many extra nodes to maintain the connectivity between adjacent sensors, the coverage-first method uses the most sensors except under grid algorithm when  $R_s < R_c < \sqrt{3} R_s$ . Also when  $R_c > \sqrt{3} R_s$ , the coverage-first method works the same as our optimal algorithm because of enough communication distance. To sum up, our optimal algorithm uses the least number of nodes to satisfy both coverage and connectivity in all four different situations.

The ratio of coverage in Region of Interest (ROI) is defined in [4] as shown in Equation (12).

$$C_r = \frac{\bigcup_{i=1, \dots, n} A_i}{A} \quad (12)$$

$A_i$  is the area covered by the  $i^{th}$  node;  $N$  is the total number of nodes;  $A$  stands for the area of the ROI, which is simulated under our optimal deployment. In our simulation, we set  $R_c = 4$ ,  $R_s = 6$  and  $R_c = 8$ ,  $R_s = 4$  to reflect the different relationship of Case 1:  $R_c \leq \sqrt{3} R_s$  and Case 2:  $R_c > \sqrt{3} R_s$  respectively. We assume that the rest of nodes exception the ones with top and mid priority die firstly. As the above-mentioned analyses because the number of nodes with top and mid priority is less than the half of all nodes in ROI, and the rest of nodes is meaningless in moving to increase the coverage, we set the maximum number of dead nodes are 120 and 150 respectively. The comparison of the coverage ratio between the network with self-healing and the other one without movement is shown in Figure 15. The red lines are the simulation in Case 1 and the blue lines represent that in Case 2. We can find our algorithm with





**Figure 15. (a) The coverage ratio of 100\*85 rectangle in two Cases, (b) The coverage ratio of particular region in two Cases.**

self-healing maintains much higher coverage ratio by contrast with the original network without movement when the nodes die continuously. The lines with our algorithm decline slowly at first and then get faster as all top priority nodes used up while the mid priority nodes start moving. When all the nodes with top and mid priority have already relocated, the slope of the line is the same as another ratio line of original network without moving. However, when most nodes have already died, the slope of original network's coverage ratio line will rise slowly as a part of top and mid priority nodes that occupy smaller area can not be available. Because of less nodes in Case 2, the coverage ratio decline faster than the ratio in Case 1.

## 6. Conclusions

In this work, the optimal deployment with self-healing movement algorithm has been proposed to ensure the coverage and connectivity of particular area by fewer sensors as compared to other three methods. Besides, with the capacity of self-healing the coverage of the entire particular area is obviously enhanced by contrast with another network without movement when some nodes are already dead. Thus, this method can be applied in closed sensing area or semi-enclosed sensing area with boundaries or obstacles which are modeled by irregular polygons or arches. Furthermore, with a combination of programs and protocols, the Ad-hoc network can be built.

## 7. References

- [1] F. Akyildiz, W. Su, Y. Sankarasubramaniam, and E. Cayirci, "A survey on sensor networks," *IEEE Communications Magazine*, Vol. 40, No. 8, pp. 102–114, August 2002.
- [2] A. Sekhar, B. S. Manoj, and C. S. R. Murthy, "Dynamic coverage maintenance algorithms for sensor networks with limited mobility," *Pervasive Computing and Communications, PerCom 2005, Third IEEE International Conference*, pp. 51–60, March 8–12, 2005.
- [3] G. Wang, G. Cao, T. La Porta, and W. Zhang, "Sensor relocation in mobile sensor networks," *Proceedings of the 24th International Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM05)* Miami, FL, March 2005.
- [4] N. Heo and P. K. Varshney, "Energy-efficient deployment of intelligent mobile sensor networks," *IEEE Transactions on Systems, Man, and Cybernetics—Part A: Systems and Humans*, Vol. 35, No. 1, pp. 78–92, 2005.
- [5] S. Shakkottai, R. Srikant, and N. Shroff, "Unreliable sensor grids: Coverage, connectivity and diameter," *Proceedings of IEEE Infocom*, San Francisco, March 2003.
- [6] S. S. Dhillon, K. Chakrabarty, and S. S. Iyengar, "Sensor placement for grid coverage under imprecise detections," *In Proceedings of the Fifth International Conference on Information Fusion*, pp. 1580–1588, 2002.
- [7] K. Kar and S. Banerjee, "Node placement for connected coverage in sensor networks," *Proceedings of the Workshop on Modeling and Optimization in Mobile, Ad Hoc and Wireless Networks (WiOpt' 03)*, Sophia Antipolis, France, 2003.
- [8] S. Ganeriwal, A. Kansal, and M. B. Srivastava, "Self aware actuation for fault repair in sensor networks," *Robotics and Automation, Proceedings, ICRA'04, IEEE International Conference*, Vol. 5, pp. 5244–5249, April 26–May 1, 2004.
- [9] J. Chen, S. Li, and Y. Sun, "Novel deployment schemes for mobile sensor networks," *Sensors*, No. 7, pp. 2907–2919, 2007.
- [10] Y. Zou and C. Krishnendu, "Sensor deployment and target localization based on virtual forces," *INFOCOM 2003, Twenty-Second Annual Joint Conference of the IEEE Computer and Communications Societies, IEEE*, Vol. 2, pp. 1293–1303, March 30–April 3, 2003.
- [11] H. N. Shu and Q. L. Liang, "Fuzzy optimization for distributed sensor deployment," *Wireless Communications and Networking Conference, IEEE*, Vol. 3, pp. 1903–1908, March 13–17, 2005.
- [12] X. L. Wu, J. S. Cho, B. J. d'Auriol, and S. Y. Lee, "Optimal deployment of mobile sensor networks and its maintenance strategy," *Computer Science, Advances in Grid and Pervasive Computing*, Springer, Berlin, pp. 112–123, June 21, 2007.
- [13] D. Pompili, T. Melodia, and I. F. Akyildiz, "Deployment analysis in underwater acoustic wireless sensor networks," *Proceedings of the ACM International Workshop on Under-Water Networks (WUWNet)*, Los Angeles, CA, September 2006.
- [14] A. Savvides, C. C. Han, and M. B. Srivastava, "Dynamic

- fine-grained localization in ad-hoc networks of sensors,” ACM MobiCom, Rome, Italy, pp. 166–179, July 2001.
- [15] S. Borbash and M. McGlynn, “Birthday protocols for low energy deployment and flexible neighbor discovery in ad hoc wireless networks,” ACM MobiHoc, Long Beach, USA, 2001.
- [16] W. B. Heinzelman, A. P. Chandrakasan, and H. Balakrishnan, “An application specific protocol architecture for wireless microsensor networks,” IEEE Transactions on Wireless Communications, Vol. 1, No. 4, pp. 660–670, 2002.
- [17] M. Younis and K. Akkaya, “Strategies and techniques for node placement in wireless sensor networks,” A Survey in ELSEVIER Ad Hoc Networks, No. 6, pp. 621–655, 2008.
- [18] Y. Zou, “Coverage-driven sensor deployment and energy-efficient information processing in wireless sensor network,” Duke University, 2004.



# An Adaptive Data Aggregation Algorithm in Wireless Sensor Network with Bursty Source

Kumar PADMANABH, Sunil Kumar VUPPALA

*Software Engineering and Technology Labs (SET Labs) Infosys Technologies Limited, Bangalore, India*

*E-mail: {Kumar\_padmanabh, sunil\_vuppala}@infosys.com*

*Received April 20, 2009; revised May 25, 2009; accepted June 30, 2009*

## Abstract

The Wireless Sensor network is distributed event based systems that differ from conventional communication network. Sensor network has severe energy constraints, redundant low data rate, and many-to-one flows. Aggregation is a technique to avoid redundant information to save energy and other resources. There are two types of aggregations. In one of the aggregation many sensor data are embedded into single packet, thus avoiding the unnecessary packet headers, this is called lossless aggregation. In the second case the sensor data goes under statistical process (average, maximum, minimum) and results are communicated to the base station, this is called lossy aggregation, because we cannot recover the original sensor data from the received aggregated packet. The number of sensor data to be aggregated in a single packet is known as degree of aggregation. The main contribution of this paper is to propose an algorithm which is adaptive to choose one of the aggregations based on scenarios and degree of aggregation based on traffic. We are also suggesting a suitable buffer management to offer best Quality of Service. Our initial experiment with NS-2 implementation shows significant energy savings by reducing the number of packets optimally at any given moment of time.

**Keywords:** Data Aggregation, Data Fusion, Congestion Control, Buffer Overflow, End to End Delay

## 1. Introduction

Wireless sensor network (WSN) is a network of sensor nodes. The main constituents of the WSN nodes are the communication devices (i.e. receiver and transmitter), a small Central Processing unit (CPU), a sensing device and a battery. The sensor node senses and gathers information from the surroundings; the CPU executes some control instructions and the communication unit sends the information to the base station through the network of such a large number of nodes.

WSN is distributed in nature and an event based system. Due to size and battery power limitations, these devices typically have limited storage capacity, limited energy resources, and limited network bandwidth. Due to these limitations, WSN differs from traditional communication networks in several ways. These limitations of sensor nodes demand specialized optimization techniques. Typically in WSN applications, a large number of Sensor Nodes (SNs) are covered over the specific target area in close proximity to each other. In such de-

ployments, spatial correlation of data is observed where neighboring sensor nodes report data values with a high degree of correlation.

Another kind of correlation observed in sensed environmental data is the temporal correlation of data where the successive sensed parameter values are found to be identical and varies slowly except in the case of unexpected events [1,2].

The spatial and temporal correlations of the WSN data can be exploited favorably for the development of efficient communication protocols in the WSN. Moreover, there is redundancy in the sensor data. The communication cost imposed due to redundant data is unnecessarily consumes lifetime of the nodes and bandwidth. In wireless sensor networks, several information can be combined together and represented by same number of bits. Once this is done the energy consumption in the communication process will be reduced. This process is known as *data aggregation*. Data aggregation schemes are the most popular way of using the correlation in sen-



sor data.

Data produced by nodes in the network propagates through other nodes in the network via wireless links. When compared to local processing of data, wireless transmission is extremely expensive. Researchers estimated that sending a single bit over radio is at least three orders of magnitude more expensive than executing a single instruction. With the new developments in the hardware of the motes, increasing memory size is giving us the chance to process the data, perform buffer management operations, so as to reduce the number of transactions over the radio.

For Scalability and flexibility of WSN applications, we need to consider this data aggregation as this results in energy saving and optimized performance. Indeed, several research efforts have been proposed in different forms of aggregation to achieve energy efficiency [1–4].

The aggregation process can be lossless or lossy. In lossless aggregation, more information is embedded into a single packet (instead of one packet for every information) thereby combining all headers into single header and same data bits. In lossy aggregation many data packets are passed through aggregation function that generates a single packet which has no information about the original data. These functions are computed by the intermediate nodes based on the data received. Thus, at each intermediate node, the amount of outgoing data is considerably lower than the amount inputted, resulting in increase of computational overhead thereby decreasing the transmitted data. The degree of aggregation (DoA) is defined as the ratio of number of bits present in all the packets considered for aggregation in one round of aggregation and the number of bits present in the aggregated packet.

There are two different types of routing in WSN literature, namely address centric and data centric. *Data centric routing* [1] is used as one of the key techniques to support in-network aggregation. Based on the data rather than the data sources and destinations, data centric routing aims to find path from multiple sources to a single destination that promote data aggregation.

Another approach is using hierarchies, where sensor nodes are usually organized into *clusters*. To perform the data aggregation nodes communicate with each other and form the clusters in order to share their sensed data. Even though such energy savings are desirable, data aggregation is sensitive with delay.

WSNs have wide range of applications. We focus on data aggregation technique that target all classes of sensor network applications from monitoring to industrial grade applications.

The rest of this paper is *organized* as follows. In Section 2, we give an overview of data aggregation techniques in WSN from the literature and motivation for our

work. We present system description and parameters in Section 3 and our approach is discussed in Section 4. Results and graphs are analyzed in Section 5 and finally the paper is concluded in Section 6.

## 2. The Related Work, Motivation and Contribution

Previous studies have proved that substantial energy savings are not only possible but essential for the success of wireless sensor networks [1]. We analyze some previous and on-going research efforts to put our work in perspective. The delay which occurred in the process of aggregation, (termed as aggregation delay) is a function of number of hops between the destination and the farthest source, and depends upon the aggregation parameter such as degree of aggregation, which will be defined in the next section. To maximize the degree of aggregation within the network, data tend to be routed through the paths that promote aggregation, rather than shortest path, which contributes additional delay.

The authors of [1] dealt with the performance issues of sensor data aggregation. They have presented a technique for delay energy trade-off in the presence of non-trivial (time consuming) aggregation. This is a mechanism to perform data centric aggregation. In their algorithm they used application specific knowledge which in turns provides a means to augmenting throughput. One of the limitations is due to its application specific approach. This algorithm is not adaptive.

The authors in [2] proposed an algorithm of aggregation which is a variant of directed diffusion. In this, intermediate nodes collect data for a specific amount of time or till they collect a fixed amount of data and send them for aggregation. The accuracy of aggregation will depend on the delay allowed at the intermediate nodes, which is specified by the application. This can improve path sharing and attain significant energy savings when the network has higher nodal density compared with the opportunistic approach. However, the idea is limited to specific amount of time or specific volume of data which is application dependent.

The authors of [3] investigate the tradeoff in the presence of *both* data aggregation and topology control (through the sleep/active dynamics of sensor nodes). In these data aggregation technologies, all aggregator nodes would wait for a fixed-period of time before performing aggregation operation. So when the time triggers, the aggregation nodes can receive responses from all of its children. This approach can save more energy consumption, but bring larger latency to the whole network.

The authors in [4] study the energy-accuracy tradeoff under two different types of aggregation: one is snapshot aggregation which is performed once, and other one is periodic aggregation which is regularly performed. The

authors claim completely distributed and localized (nodes exchange information only with immediate one hop neighbors) algorithm, however the parent should receive an exact number of messages, equal to the number of its children and the final result is only available at the user node. Snapshot aggregation on the other hand is very sensitive to the stability of the hierarchical structure.

The work by the authors in [5] provided a new stochastic decision framework to study the fundamental energy-delay tradeoff in distributed data aggregation. Adaptive real-time dynamic programming (ARTDP) is asynchronous value iteration scheme and is suitable for on-line implementation only. This scheme might be good to have energy-delay trade-off case but Adaptive Application-Independent Data Aggregation (AIDA) [6] offer better energy benefits than this scheme. The authors of paper [6] describe an aggregation scheme that adaptively performs application independent data aggregation in a time sensitive manner. AIDA performs *lossless* aggregation by concatenating network units into larger payloads that are sent to the MAC layer for transmission. This may not suite all the applications.

Some aggregate functions require the concatenation of all readings to be returned to the host node. For example, in order to accurately determine the median value in a network [7], the host node must know all the values. In this case, it may still be possible to reduce the size and number of messages by applying compression. Researchers propose a unique data structure called a Quantile Digest (q-digest), which provides approximate results that adhere to a strict error bound. But it is a good approximation scheme when there are wide variations in frequencies of different values.

The work by authors of [8] handles the case of lossy aggregation while bounding the number of messages transmitted in the network. They propose a Marginal Gains Adjustment (MGA) algorithm for the problem of bandwidth constrained aggregate continuous queries over sensor network. This does not consider all cases of aggregation and is not adaptive in nature.

Sometimes application specific aggregation will be giving better results rather than the general schemes as it can understand the environment conditions better. So we need to consider some application knowledge and propose a general purpose aggregation scheme.

## 2.1. Motivation

Even though several research works in the literature have discussed the problems and approaches of developing data aggregation processes mainly for energy, bandwidth and memory space savings by minimizing the data transferred in sensor networks [1,2]), however authors of these papers fail to address following practical problems:

Quality of Service (QoS) issues in data aggregation: In

sensor network there are several types of data. Namely normal hello packets, normal sensor data packets, some important alert data packets and control messages from the base station. The control message from the base station and the alert sensitive data packets are very important in nature and QoS provided to these packets should be better than others.

Adaptive mechanism: The parameter of the data aggregation such as DoA, QoS cannot be decided and fixed due to the burst nature of the sensor network. It should be adaptable enough. Feedback should be there to make the system controlled and adaptable. Though there is some paper available but they don't address QoS and adaptable aggregation simultaneously.

Scheduling: In the process of addressing QoS, we need to schedule the packets and apply some of the buffer management policies before applying aggregation process.

Most of the proposals in the literature give modeling and simulation of the WSN scenario for various parameters like energy, priority, delay, degree of aggregation supported with the mathematical proofs. The authors of these papers have considered either distinct parameter in each piece of work separately or they have considered only few parameters together [1,2].

Moreover these proposed methods are too complex to be implemented in hardware of current state of the art. Although several schemes for programming and data aggregation in WSNs have been proposed in literature, few actually provide experimental validation and performance evaluation [5,6].

So there is a need to design a data aggregation mechanism in WSN by considering different QoS parameters and take the feedback mechanism to make the system adaptive and save the energy. This general purpose data aggregation should be able to apply for all WSN applications, considering the priority information and application knowledge for aggregation function.

Our approach is to have buffer management in the aggregator nodes to make the adaptive algorithm obeys the rule that degree of aggregation is proportional to number of packets. Special packet formats are considered in the aggregation. So this approach can be used for wide range of sensor applications.

## 2.2. Contribution of This Paper

There are considerable amount of work in data aggregation available in existing literature. Authors of these papers dealt with application dependent or adaptable technique, QoS issues, related techniques in separately. In this scenario following is the contribution of this paper:

1) Lossy Lossless Aggregation: In the same algorithm lossy and lossless aggregation has been taken care. Depending upon the requirement algorithm switched from

lossy to lossless.

2) Controlled Degree of Aggregation: The degree of aggregation is control parameter and existing number of packets in the buffer determine the instantaneous value of degree of aggregation.

3) Buffer Management: In the same algorithm we have taken care of buffer management which optimizes the QoS by minimizing the packet loss due to buffer overflow.

The above three points are our unique contribution in this.

### 3. System Description

We consider two types of nodes in our system, normal nodes and aggregating nodes. Normal nodes do not perform aggregation. They sense the data and send it to the sink. They also forward the data generated by other nodes. Aggregating nodes work as normal nodes and perform aggregation. Only local aggregation can be done at normal nodes.

An aggregator receives the data from one or more normal nodes, performs an aggregation based on the algorithm and then forwards the aggregated packet. In WSN, data from all of the nodes are supposed to be shipped to the base station only. Thus a base station in the WSN is a typical sink where the data reaches finally. Actually this base station connects the individual sensor node to outside world.

In our system we consider following four types of packets:

Hello packets which consists of the information about the source nodes and may contain the routing information. It does not contain any sensor data or any other data.

The control packets contain some of the control parameters. It may originate from the base station or from other nodes. The control parameter may be some system control instruction or to set some flag or otherwise.

Normal data packets: In the sensor network the data packets are formed with sensor reading and headers. Regular messages are those messages that contain such sensor data which fall in the expected range. Typically it is the instantaneous sensor reading. In this case sensing is being done as a regular practice which occurs without any event of interest.

Critical data packets: Critical data packets are those packets which contains sensor data and header. This sensor data is generated with an event of interest. For example, the normal temperature of office workplace is 25°C. A packet with sensor reading of 25°C will be known as normal packet. However if the sensor reads a temperature of 75°C it will be an event of interest and the packet which contains this reading will be termed as critical packet.

We assume whether a particular node will work as a normal node or aggregator nodes is decided by some technique which is not in the scope of this work.

Typically there are two way of aggregation. Firstly extracting the sensor data from the packet and considering many such sensor data to pass through an aggregation function to get a single data. For example if the sensor data is temperature reading then we can consider many temperature readings to take average of them. Thus, before aggregation we have multiple sensor data however after aggregation we have a single average value. When this average value is used to form a packet we call it as an aggregated packet. This aggregated packet is lossy. Because at receiving end we cannot reproduce the original sensor data with the average value of reading. Therefore it is called as lossy aggregation. However there is another technique in which sensor readings are extracted from multiple packets and they are put into single packet with one header only. Here nothing is lost however we are getting rid of header information. The packet length will be variable in this case. This is lossless aggregation.

In our system we consider both lossy and lossless types of aggregation. To choose between lossy and lossless is completely application dependent. However general rule is that when packet size is not fixed, we can go for lossless aggregation and when we have an optimally designed fixed size packet we can go for lossy aggregation.

In our system we have considered two different level of aggregation taking place at different nodes. It starts from the source node itself. The first among these two levels of aggregation is local aggregation. Here any particular node generates data from sensor readings and put them into a packet. Nodes may decide to put more than one sensor data in single packet; they may take average, min-max of some of the data actually depending upon the application and then put them into a single packet. So the number of data packets is reduced and information of many possible data packets is embedded into a single data packet. We call it as local aggregation or *level-1 aggregation*. It is to be noted that though it is a local aggregation, this is a global policy of data aggregation. It means all other nodes of similar kind will do same aggregation throughout the network.

Second level of aggregation happens with the data packets of locally aggregated data down the line towards the base station. Thus it may happen at any intermediate node from the source node to the base station. In this second level of aggregation some aggregation function is applied to the data streaming from various source nodes to these level-2 aggregator nodes for lossy aggregation or sensor data are extracted from the packets to put them into single packet for lossless aggregation. This again depends upon the application. In this case aggregation

may be done on one kind of sensor data.

### 3.1. Useful Parameters Considered in the System

The performance evaluation of the data aggregation mechanism can be done by analyzing some of the parameters. In our system we consider following parameters:

#### 3.1.1. Degree of Aggregation

We define degree of aggregation as a ratio of total number of number of bits in all packets considered for one round of aggregation process and total number of bits in aggregated packets.

Let us consider that  $X$  is the number of data bits in the packet and  $H$  is the number of header bits in a single packet. Thus if  $n$  number of packets are considered for aggregation in one round of algorithm of aggregation, let us consider the lossy and lossless aggregation case separately to define degree of aggregation formally,

Lossy aggregation: In this case  $n$  numbers of sensor data are passed through aggregation function to get a single packet. Additionally  $z_1$  number of additional bits will be added to form the aggregated bits.  $z_1$  is the number of bits required to carry the aggregation information like average value, statistical value, number of packets involved in the aggregation. This  $z_1$  can be fixed for a particular application. These  $z_1$  bits are used to decode the aggregated sensor data at the sink. Therefore DoA in this case will be defined as:  $DoA = \frac{n(X+H)}{X+H+z_1}$ , for a very minimal additional bits

as an identifier i.e.  $z_1 \ll n(X+H)$ . The degree of aggregation will be reduced to  $n$  itself.

In the case of lossless data aggregation the degree of aggregation is defined as similarly. However the number of bits after aggregation will be reduced to  $nX + H + z_2$ . Therefore degree of aggregation for lossless aggregation can be defined as  $DoA = \frac{n(X+H)}{nX+H+z_2}$  it is to be noted

that degree of aggregation for lossless case is lesser than its lossy counterpart.

#### 3.1.2. QoS

We have considered priority based service to four types of packets defined earlier. We consider hello packet and normal data packets as general packets. Other packets, namely, control packets and critical packets are important packets. The important packets will have priority over the normal packets for service. We apply a special buffer management policy with data aggregation to achieve this. Typically we don't want to have lossy ag-

gregation or loose any packets from these high priority packets.

#### 3.1.3. Packet Format

To achieve the QoS discussed earlier we have proposed a general packet format which is applicable for both lossy as well as lossless aggregation. In WSN, there is no fixed format for the packet in practice. We are proposing both fixed and variable length packet format.

#### 3.1.4. Fixed Packet Format

For lossy aggregation following is the packet format considered. We have a typical OS based header packet type and data field. It is to be noted that data packet have fixed length in this case.

For example, TinyOS [9] default payload is of 29 bytes. TinyOS Header field consists of destination address, type, group id and message length. Rest of the payload is defaulted to 29 bytes. In our packet structure of multi hop routing, along with standard TinyOS header, we have few more fields as additional header, namely source node address, parent node address, hop count, sequence number and last forwarder id. Rest of the payload consists of different sensor analog to digital converter (ADC) values indicating sensor data readings. The packet is represented in associated Figure 1.

As the payload is taken as fixed size for the aggregated packet in lossy aggregation, one extra type field is enough to differentiate normal packet and aggregated packet.

#### 3.1.5. Variable Length Packet Format

We propose special adaptable packet format here. The header field will be the same except there will be additional fields in header which will carry information about the length of the packet. The length of data field will be variable so the total length of the packet will be variable in nature and adapt to the current scenario. This is mainly

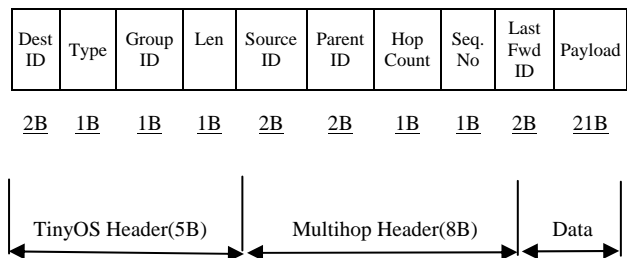


Figure 1. Packet structure for TinyOS with multi hop routing.

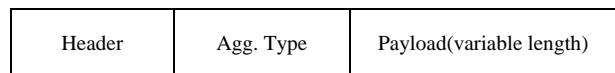


Figure 2. Adaptive aggregated packet.

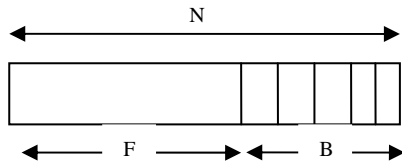


Figure 3. Buffer and DoA relationship.

used for lossless aggregation. However it is to be noted that there is a maximum limit to the length of this payload. For example, in case of TinyOS, the message length can be up to 116 bytes. So there will be different combinations possible to prepare the variable aggregated packet as sensor readings from different nodes need to be sent in a single packet. We shall consider both the type field in the header and the aggregation type in payload to handle different combinations (Figure 2). This variable packet needs to be interpreted correctly at the base station by considering the aggregator type filed. By this way, system can generate aggregated packet on the fly based on the inputs given to the system.

### 3.1.6. Special Buffer Management for Data Aggregation

Data aggregation involves combining several sensor readings in intermediate nodes. This in turn requires storing the packets from different sensor nodes and processes them in the memory space available in nodes and outputting aggregated packet. To input the packets from different sensor nodes, we consider buffer space in the aggregator node and to process them we need a special management policy [10,11] so that it can provide specific number of packets to aggregation process after considering type of packet and DoA type.

Buffer acts as a storage for the packets and works similar to a queue. In our system, we consider a temporary buffer and multiple queue system in main buffer. First the input packet reaches the temporary buffer and then caters to different priority queues. We define different queues for different priority packets. For the first queue in the buffer, we push normal and Type-1 critical packets for which aggregation is needed. Second queue is for important packets and third queue is for critical packets. Let us consider that  $N$  is the total space in buffer and  $B$  is the number of packets in the buffer (Figure 3). Thus,  $F$  is considered as the difference between  $N$  and  $B$  (i.e.  $N-B$ ) indicating free space in the buffer [10].

The policy considered in the buffer management follows these rules.

- 1) General packet processing is on the first come first serve basis.
- 2) From temporary buffer, the packet is pushed to relevant queue in the main buffer based on the type of packet.
- 3) A packet is never dropped as long as there is room

in the main buffer.

4) A packet from temporary buffer is discarded only if the main buffer is full.

5) DoA is proportional to the number of packets ( $\text{DoA} \propto B$ ) as shown in the Figure 3.

6) If there is no space for incoming packet, packet of the low priority is dropped from the temporary buffer.

## 4. Our Approach

In this section, we present our approach of adaptive data aggregation based on different parameters mentioned earlier in the paper. In our system, aggregation is performed in two levels after storing and processing the sensor data packets in the buffer.

Hello packets and control packets are processed without aggregation. Aggregation is performed for the normal packets. Based on the application demand, critical data packets can also be aggregated. If the node can take necessary action in response to the event of interest, we may send the critical data packet after the aggregation, referred as Type-1 critical packets. This is implemented by incorporating necessary functionality inside the node. For these Type-1 critical packets, a control packet is also generated from the node. This control packet could be sending an alarm signal or sending an alert to the corresponding person. For Type-2 critical packets, no aggregation takes place as the critical data packet is sent to the sink as soon as possible. In this case, the sink responds to the received critical data packet, which is generated for the event of interest from the node.

The sensors sense the data at frequent intervals of time and check for the possibility of any local aggregation before generating the packets. This is referred to as local aggregation. After local aggregation, the packet is generated and enters into the buffer of the next node towards the sink from the input queue. Based on the aggregation mechanism and type of packet, few packets are processed and an aggregated packet is outputted as described in the algorithm. The effectiveness of data aggregation is improved by taking feedback from the system. This feedback contains the number of packets to be considered for aggregation in each round. We consider the feedback and degree of aggregation type in the buffer management to make adaptive aggregation as shown in Figure 4.

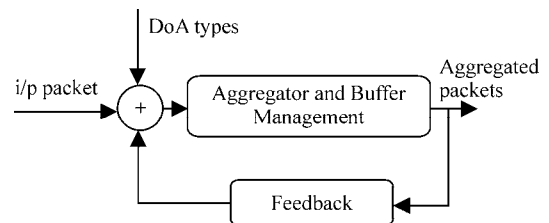


Figure 4. Feedback mechanism in the data aggregation.

Range of packets	DoA Type
$B < M1$	1
$M1 < B < M2$	2
$M2 < B < M3$	3
$M3 < B < N$	4

**Figure 5. Degree of aggregation type based on count/time and number of packets.**

If aggregation is being done in a control environment, then degree of aggregation will not be a fixed parameter in the system. It will be adaptive to the instantaneous requirement of application. Moreover, the system can not aggregate all the packets present in the buffer due to the processing involved, which is delay sensitive. So this leads to requirement of different DoA types. The system can be either packet-count based or time based. The system waits until the buffer reaches the specified number of packets in the count based type, where as in the time based type the system waits for a particular amount of time which in turn decides automatically the value of DoA types from Figure 5.

We choose the number of packets to be aggregated at each instance and are given to the system as feedback so that corresponding DoA type is chosen to decide the aggregation mechanism to be performed. For example, let  $M1$ ,  $M2$ ,  $M3$  represent different numbers, which are selected based on the application. If the present number of packets in the buffer is less than  $M1$  choose DoA Type-1. DoA Type-2 is chosen if the number of packets lies between  $M1$  and  $M2$  as shown in Figure 5.

The DoA type and the range of packets can be adaptive based on the feedback from the system so that we can optimize the aggregation output.

For each round of operation, specified number of packets are aggregated based on the above mentioned considerations as described in the algorithm. The resulting aggregated packet is sent as output from the system. This holds good for both lossy and lossless aggregation. From the aggregated packet, we calculate the DoA based on the number of packets involved in the aggregation and type of aggregation like lossy or lossless which is explained in the previous section.

#### 4.1. The Algorithm

By considering all the parameters and features mentioned in the last section, we propose an adaptive algorithm for data aggregation in two levels for both lossy and lossless types of aggregation. Level-1 aggregation is being performed locally just after reading the sensor data. However, level-2 aggregation is being performed on the sensor data coming from various nodes. We follow the algo-

rithm for level-1 and level-2 aggregation as explained in Table 1. In level-2 aggregation, we logically divide the system into two phases namely, collection and aggregation phases. Collection phase collects the data to be aggregated where as aggregation phase processes the actual aggregation. The collection and aggregation phase repeat until the system is running. Few steps will be common for both lossy and lossless aggregation.

Let us first consider the level-1 aggregation. The sensor nodes sense the data from the environment at frequent intervals of time. After sensing the data, it checks

**Table 1. Algorithm for data aggregation.**

```

Level-1 Aggregation:
Require: Sensed data
{
    if (local aggregation) then
        if (event of interest) then
            Generate packet; Forward packet to Sink
        else
            Store sensed data and aggregate with
next readings
        end if
    else
        Generate packet
    end if
    Forward to Aggregator
} //end level-1

Level-2 aggregation:
Collection Phase:
Require: packet reaches aggregator
{
    Store the packets in buffer
    if (packet priority = Critical/Important) then
        Forward packet to sink (no aggregation)
    else
        Wait for T Sec/Count M.
        if (Time/Count reached) then
            Apply aggregation
        end if
    end if
} //end collection phase

Aggregation phase:
Require: Number of packets and DoA Type
{
    Take the number of packets to aggregate
    (Feedback parameter in next iterations)
    Extract the sensor data from different packets
    if (lossless aggregation) then
        Format the packets with new type
        Aggregate the packet and send to sink; Compute
DoA in bits
    end if
    if (lossy aggregation) then
        Use aggregation function; Compute DoA in bits
        if (aggregation function = Min/Max Type)
            Continue the aggregation in next hops
        until packet reaches sink
    else
        Send the aggregated packet to Sink
    end if
end if
}

```

for the need of any local aggregation (level-1) within the node. If local aggregation is possible, it stores the sensor data and waits for the next sensor readings before generating the packet. In case of any event of interest, the packet is generated without waiting for local aggregation and forwarded to the next node towards sink. Otherwise, this node generates the packet and forwards it to the aggregator node. If local aggregation is performed, it is indicated by the type field in the packet format. So at this point, several packets from different nodes reach the aggregator nodes. This is referred as level-1 aggregation.

In the collection phase of level-2 aggregation, the aggregator node collects the packets in the buffer. It checks for the priority of the packets and append the packet in the buffer as per the priority described in earlier section. If the packet is found to be critical or important, they are forwarded from the aggregator towards sink with out any aggregation. In other words those packets are not aggregated at all. In the system, it needs to identify either to follow count based or time based mechanism for the aggregation. In the count based, the system will wait for specific number of incoming packets to be inserted into the buffer. Then it aggregates the fixed number of packets (as per current DoA) from the head of the queue. The choice of count based or time based depends on the application.

Actually in phase-2, DoA type and number of packets to be aggregated are taken as inputs. DoA type is taken from predefined readings of the system as given in Figure 5. The number of packets to be aggregated at any particular moment of time is determined by current space in buffer which is taken as a feedback parameter as shown in Figure 4. All these steps are similar for lossless and lossy both. However, from this point onwards, lossy and lossless aggregation methods differ and are described as follows:

In lossy aggregation, particular number of packets in buffer is considered for aggregation from the collection phase and the sensor readings are extracted from different packets. Then according to requirement of application, a particular aggregation function is selected.

Basically there are two types of aggregation functions possible. Functions like average, standard deviation are limited to one hop only in aggregation process. That means, once the packets are aggregated with this function, no further aggregation is suggested till it reaches the sink. In the other case, functions like minimum, maximum can continue aggregation till it reaches to sink, further reducing the number of packets transferred in the system.

It is to be noted that except the type field for indication of aggregated packet, the size of the packet remains same in this case. At each aggregation step, the DoA is computed in terms of bits. In our system we have considered this for energy saving calculation and to analyze the system performance.

In case of lossless aggregation, the sensor readings are extracted from the packets taken from collection phase. All these sensor readings are aggregated and formatted into a packet with a new type and variable length. The type and length fields describe the packet format to retrieve the readings at the sink. The DoA is computed in terms of bits, at each aggregation step.

Once the packets with or without aggregation reach the sink, it extracts the readings based on the packet type and is used for the application. These steps are described in the algorithm shown in Figure 6 and Table 1.

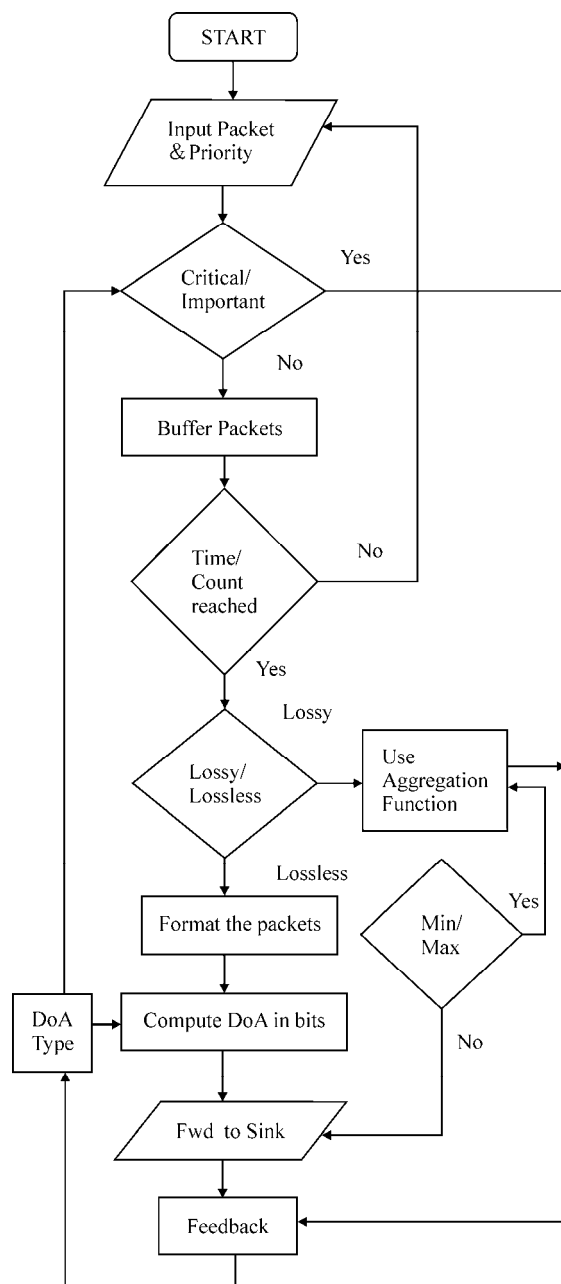


Figure 6. Flow chart of level-2 aggregation.

Different applications can be taken up to illustrate our scheme. Let us consider temperature monitoring application with average value as the aggregation function as one illustrative application. This is a typical monitoring application using sensor nodes. Nodes are deployed in a room or in an open space where there is a need for monitoring. The sink or base station is located either in the same room where the nodes are deployed or in another room. The packets carrying the sensor data reach the base station in regular intervals of time. The base station process the data and business logic is applied. So the critical data packets should be sent to the base station as soon as possible without the aggregation. Required alerts are raised to the concerned person if the temperature readings are out of bound. Out of bound temperature readings are considered as the event of interest. As the aggregation function considered is average, the aggregation mechanism is considered up to 1 hop level only. After the packets are aggregated, they are sent to sink without further aggregation in the next levels.

In this application, the packet is generated every one second at each node indicating the temperature reading. If there is not much change in the sensor reading from the previous value (up to a reference) we can do local aggregation. After that the packets are generated and reach the aggregator.

Based on the algorithm of lossy aggregation, packets are aggregated and the aggregated packet is indicated as a special type of packet. We follow the table (Figure 7) to choose the DoA type.

Here it is a time based function. For every 10 seconds the aggregation algorithm is called to check the number of packets (B) and DoA type to apply the aggregation mechanism in the system.

## 5. Analysis & Results

In this section we analyze the results from simulation of our model. Different parameters considered in the system are defined as follows:

**Average Delay:** Delay is taken as the time each packet is in the buffer in the process of aggregation. Average delay is calculated taking average time each packet spends in the buffer.

**Degree of Aggregation:** The DoA is defined as the ratio of number of bits present in all the packets considered for aggregation in one round of aggregation and the number of bits present in aggregated packet.

**Packet Loss:** It indicates the number of packet drops or loss due to buffering of the packets to aggregate in the process of aggregation. Critical packets, important packets and normal packets are treated differently in the buffer and corresponding loss rate is considered.

Range of packets	Time	DoA Type
$B < 10$	10 sec	1
$10 < B < 20$	10 sec	2
$20 < B < 30$	10 sec	3
$30 < B < 100$	10 sec	4

Figure 7. Example of time based DoA.

These parameters are considered for different Constant Bit Rate traffic (CBR) traffic, network sizes in both lossy, lossless aggregation based on count and time. We have conducted simulations in Network Simulator (NS-2) [12] to test the performance of our model in large scale. Packet size is taken as 32 bytes, as described in earlier section. Lossless aggregated packet size is variable and maximum size is considered as 116 bytes. So each lossless aggregated packet can accommodate maximum of 20 packets considering 3 sensor readings per each packet. CBR varies from 1 packet/sec to 40 packets/sec. Different network sizes from 10 sensor nodes to 500 sensor nodes are considered in the simulation. Buffer can accommodate a total of 300 packets. For count based, 10 packets are aggregated at a time. In time based, we have taken interval of 1 minute, so the number of packets differs as traffic and network size increases. This buffer is divided into 3 equal parts (100 packets) for critical, important and normal packets. But this memory is sharable among these packets giving the order of preferences, as described in the system description of the paper.

### 5.1. Degree of Aggregation for Different Network Sizes

Here DoA is considered for different network sizes as shown in Figure 8 keeping the CBR as constant and tested for all four possible combinations of lossless count based, lossless time based, lossy count based and lossy

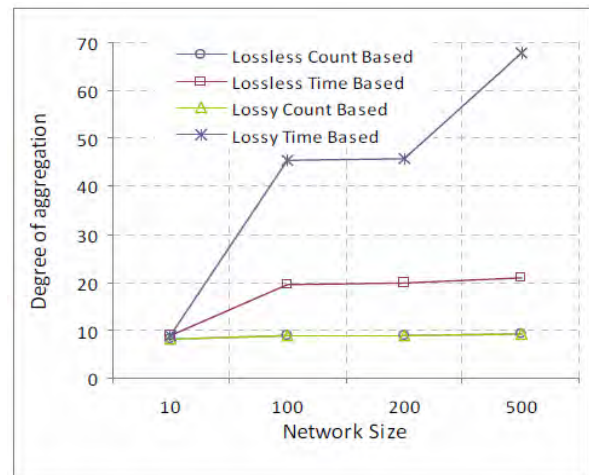


Figure 8. Degree of aggregation for different network sizes.



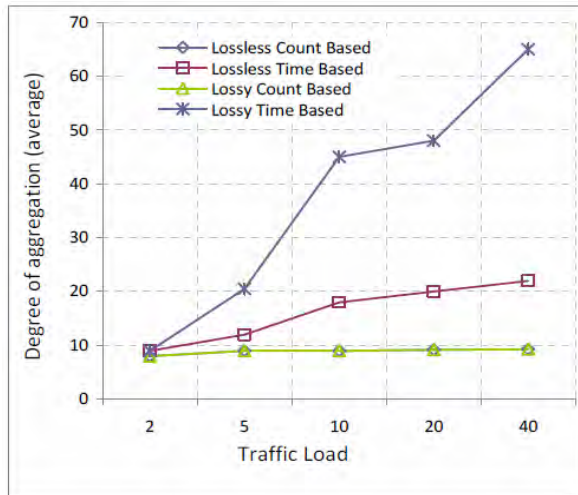


Figure 9. Degree of aggregation for different traffic.

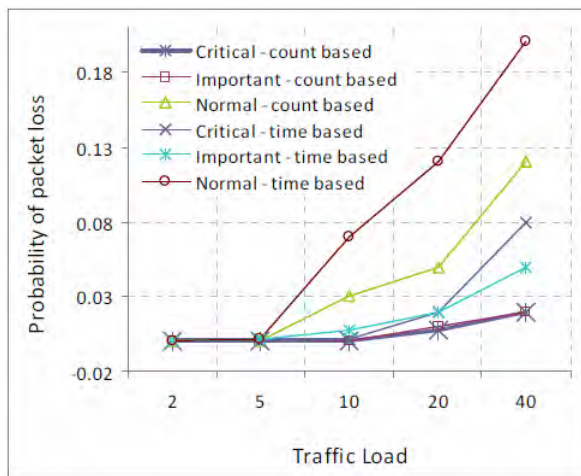


Figure 10. Packet loss for different traffic.

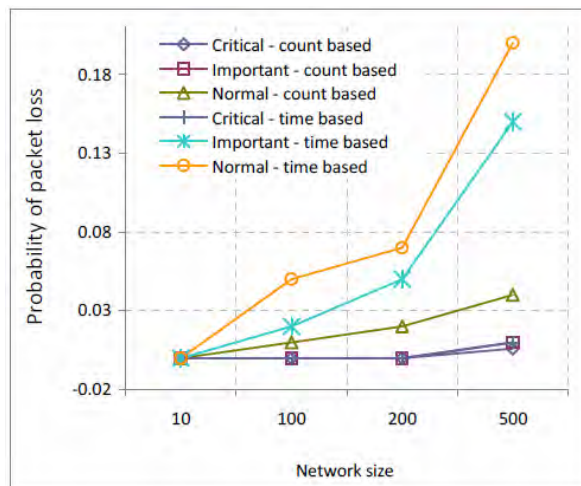


Figure 11. Packet loss for different network sizes.

time based aggregations. In count based aggregation, the DoA is in the range of 10 only both for lossy and lossless due to the fact that as soon as count reaches for 10 packets, aggregation is applied. In case of time based aggregation, DoA increases as network size increases based on the table (Figure 5) for different DoA types. In case of lossless aggregation, there is a limit of DoA as it can accommodate maximum of 20 packets, so limiting the DoA around 20. In the case of lossy aggregation, the DoA grows as per the DoA type (Figure 5) and is limited by the buffer size only. It is very evident that our proposed algorithm makes the system adaptable to instantaneous condition and the required aggregated packet is generated with proper DoA.

## 5.2. Degree of Aggregation for Different Traffic

Here DoA is considered for different traffic rates by varying CBR flow as shown in Figure 9 keeping the network size as constant of 100 nodes and tested for all four possible combinations. In count based aggregation, the DoA is in the range of 10 only both for lossy and lossless due to the fact that count of 10 is the limit to trigger the aggregation process. But DoA increases as traffic rate increases in case of time based aggregation. In time based lossless aggregation, DoA is around 20 as described in the first graph. For lossy aggregation, DoA increases as traffic load increases and grows as per the DoA type. In this case, more packets are available for aggregation with increase in traffic load. Only limitation for DoA is the buffer size. Feedback is used at each stage to choose the specific DoA type as mentioned in the algorithm.

## 5.3. Packet Loss for Different Traffic and Network Sizes

In our system, we have different types of packets (critical, important and normal) which are treated differently inside the buffer in the process of aggregation. Our goal is to minimize the loss of packets in the buffer and to have less delay, for which a trade off is required. In Figure 10, packet loss is shown for different traffic load keeping the network size as constant at 100 nodes. In Figure 11, packet loss is shown for different network sizes by keeping the CBR as constant of 10 packets per sec. In both the cases, loss of critical packets is very less initially but as network grows there are a bit of drop in the critical packets. In case of the important packets and normal packets also as traffic rate/network size increases, there is a drop in the packets due to the limitation of buffer size. But packet loss is more in time based when compared to count based as the packets in the buffer are limited in case of count based aggregation mechanism. In count based, aggregation process is triggered by reaching

a specific number of packets in buffer even more packets are generated in the system. In case of time based more packets reach buffer as traffic load increases, so the loss of packets. Overall, critical packets loss is very minimal, as desired.

#### 5.4. Average Delay for Different Network Sizes

Aggregation is applied more frequently in count based therefore least delay is observed. In time based aggregation, more number of packets gets accumulated which is resulted in more delay. However in lossy aggregation, since the number of bits reduced drastically therefore, least delay is there. In lossless aggregation no bits of information is lost therefore more queuing delay is introduced. This is why time based lossless aggregation has highest delay. All four cases can be easily interpreted from Figure 12.

We need a tradeoff between the delay and the DoA keeping the track of packet loss. Even though count based aggregation is having less delay and less packet loss rate, DoA is limited. But in case of time based aggregation, we have better DoA at the cost of delay and some packet loss. So lossy or lossless, count or time based aggregation is to be chosen based on the application type. Adaptive algorithm takes major role in getting the increased DoA based on the feedback.

#### 6. Conclusions

The idea of the aggregation is to aggregate the data required close to the source or at intermediate nodes on the way to sink instead of sending all the sensor readings through the network. In this paper the idea of lossy and lossless aggregation has been proposed within a single algorithm. The proposed algorithm works for both lossy and lossless aggregation depending upon the requirement

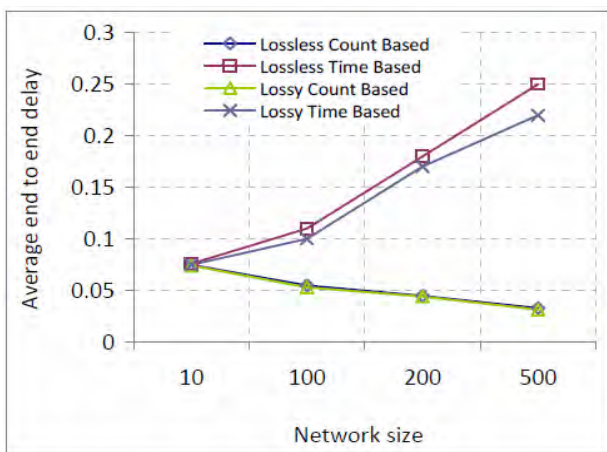


Figure 12. Average delay in aggregation for different network sizes.

and makes the system adaptive to changes which can be adjusted with the load in the buffer and buffer management policy. The ultimate aim is to offer best QoS and significant savings in the energy and number of packets to be transmitted. The experiment has been carried out with Network simulator for large scale sensor network which advocates our proposed algorithm.

#### 7. References

- [1] B. Krishnamachari, D. Estrin, and S. Wicker, "Modelling data-centric routing in wireless sensor networks," in Proceedings of the IEEE INFOCOM, 2002.
- [2] C. Intanagonwiwat, D. Estrin, R. Govindan, and J. Heidemann, "Impact of network density on data aggregation in wireless sensor networks," in Proceedings of the 22nd International Conference on Distributed Computing Systems (ICDCS'02), July 2002.
- [3] V. Erramilli, I. Matta, and A. Bestavros, "On the interaction between data aggregation and topology control in wireless sensor networks," First Annual IEEE Communications Society Conference on Sensor and Ad Hoc Communications and Networks, IEEE SECON, pp. 557–565, 2004.
- [4] A. Boulis, S. Ganeriwal, and M. B. Srivastava, "Aggregation in sensor networks: An energy-accuracy trade-off," First IEEE International Workshop Sensor Network Protocols and Applications (SNPA'03), May 2003.
- [5] Z. Ye, A. A. Abouzeid, and J. Ai, "Optimal policies for distributed data aggregation in wireless sensor networks," in Proceedings of 26th Annual IEEE Conference on Computer Communications (INFOCOM'07), Anchorage, Alaska, USA, May 6–12, 2007.
- [6] T. He, B. Blum, J. Stankovic, and T. Abdelzaher, "AIDA: Adaptive application independent data aggregation in wireless sensor networks," ACM Transaction on Sensor Network, 2003.
- [7] N. Shrivastava, C. Buragohain, D. Agrawal, and S. Suri, "Medians and beyond: New aggregation techniques for sensor networks," in Proceedings of the Second ACM Conference on Embedded Networked Sensor Systems (SenSys'04), August 2004.
- [8] A. Deligiannakis, Y. Kotidis, and Roussopoulos, "Bandwidth-constrained queries in sensor networks," The VLDB Journal, Vol. 17, No. 3, pp. 443–467, 2008.
- [9] TinyOS, <http://www.tinyos.net/>.
- [10] K. Padmanabh and R. Roy, "Cost sensitive pushout policy and expelling policies with dynamic threshold for the buffer management in differentiated service switch for versatile traffic," IEEE International Conference on Networking, Mauritius, April 23–28, 2006.
- [11] G. J. Fosdini and B. Gopinath, "Sharing memory optimally," IEEE Transaction on Communication, Vol. 31, No. 3, pp. 352–360, March 1983.
- [12] Network Simulator 2, <http://www.isi.edu/nsnam/ns/>.



The 6<sup>th</sup> International Conference on Wireless  
Communications, Networking and Mobile Computing

**September 23-25, 2010, Chengdu, China**

<http://www.wicom-meeting.org/2010>

## Call for Papers

**WiCOM** serves as a forum for wireless communications researchers, industry professionals, and academics interested in the latest development and design of wireless systems. In 2010, WiCOM will be held in **Chengdu**, China. You are invited to submit papers in all areas of wireless communications, networking, mobile computing and applications. All papers accepted will be included in IEEE Xplore and indexed by EI Compendex and ISTP.

### Wireless Communications

- B3G and 4G Technologies
- MIMO and OFDM
- Cognitive Radio
- Coding, Detection and Modulation
- Signal Processing
- Channel Model and Characterization
- Antenna and Circuit

### Network Technologies

- Ad hoc and Mesh Networks
- Wireless Sensor Networks

- RFID, Bluetooth and 802.1x Technologies
- Network Protocol and Congestion Control
- QoS and Traffic Analysis
- Network Security
- Multimedia in Wireless Networks

### Services and Application

- Applications and Value-Added Services
- Location Based Services
- Authentication, Authorization and Billing
- Data Management
- Mobile Computing Systems

### Submission Requirement:

The working language of the conference is English. All the papers must be submitted in IEEE electronic format. Instructions and full information on the conference are posted on the conference website. Anyone wishing to propose a special session or a tutorial should contact us: [wicom@scirp.org](mailto:wicom@scirp.org)

### Important Dates:

Paper Due: Feb. 28, 2010

Acceptance Notification: May. 4, 2010

### Contact Information:

Website: <http://www.wicom-meeting.org/2010>

E-mail: [wicom@scirp.org](mailto:wicom@scirp.org)



# Wireless Sensor Network (WSN)

## *Call For Papers*

<http://www.scirp.org/journal/wsn>

ISSN 1945-3078 (Print) ISSN 1945-3086 (Online)

WSN is an international refereed journal dedicated to the latest advancement of wireless sensor network and applications. The goal of this journal is to keep a record of the state-of-the-art research and promote the research work in these areas.

### **Editor-in-Chief**

Dr. Kosai Raoof , GIPSA LAB, University of Joseph Fourier, Grenoble, France

### **Subject Coverage**

This journal invites original research and review papers that address the following issues in wireless sensor networks. Topics of interest are (but not limited to):

- Network Architecture and Protocols
- Self-Organization and Synchronization
- Quality of Service
- Data Processing, Storage and Management
- Network Planning, Provisioning and Deployment
- Integration with Other System
- Software Platforms and Development Tools
- Routing and Data Dissemination
- Energy Conservation and Management
- Security and Privacy
- Developments and Applications
- Network Simulation and Platforms

We are also interested in short papers (letters) that clearly address a specific problem, and short survey or position papers that sketch the results or problems on a specific topic. Authors of selected short papers would be invited to write a regular paper on the same topic for future issues of the WSN.

### **Notes for Intending Authors**

Submitted papers should not have been previously published nor be currently under consideration for publication elsewhere. Paper submission will be handled electronically through the website. All papers are refereed through a peer review process. Authors are responsible for having their papers checked for style and grammar prior to submission to WSN. Papers may be rejected if the language is not satisfactory. For more details about the submissions, please access the website.

### **Website and E-Mail**

<http://www.scirp.org/journal/wsn>

Email: [wsn@scirp.org](mailto:wsn@scirp.org)



## TABLE OF CONTENTS

**Volume 1    Number 3**

**October 2009**

**A Cognitive Radio Receiver Supporting Wide-Band Sensing**

V. Blaschke, T. Renk, F. K. Jondral..... 123

**A Caching Scheme for Session Setup in IMS Network**

Y. F. Cao, J. X. Liao, Q. Qi, X. M. Zhu..... 132

**Metrics and Algorithms for Scheduling of Data Dissemination in Mesh Units**

**Assisted Vehicular Networks**

Z. Y. LIU, B. LIU, W. YAN..... 142

**High Resolution MIMO-HFSWR Radar Using Sparse Frequency Waveforms**

G. H. Wang, Y. L. Lu..... 152

**ContSteg: Contourlet-Based Steganography Method**

H. Sajedi, M. Jamzad..... 163

**Research on DOA Estimation of Multi-Component LFM Signals**

**Based on the FRFT**

H. T. Qu, R. H. Wang, W. Qu, P. Zhao..... 171

**Novel Rate-Control Algorithm Based on TM5 Framework**

Z. J. ZHU, Y. Q. BAI, Z. Y. DUAN, F. LIANG..... 182

**Generation of Multiple Weights in the Opportunistic Beamforming Systems**

G. Y. Lu, L. Zhang, H. Q. Yu, C. Shao..... 189

**The Effect of Notch Filter on RFI Suppression**

W. G. Chang, J. Y. Li, X. Y. Li..... 196

**Reconfigure ZigBee Network Based on System Design**

Y. Xu, S. B. Qiu, M. Hou..... 206

**Optimal Deployment with Self-Healing Movement Algorithm for Particular**

**Region in Wireless Sensor Network**

F. Zhu, H. L. Liu, S. G. Liu, J. Zhan..... 212

**An Adaptive Data Aggregation Algorithm in Wireless Sensor**

**Network with Bursty Source**

K. Padmanabh, S. K. Vuppala..... 222

